

反辞書の 効率的な構築手法に 関する研究

Study on Efficient Algorithms for Construction of
an Antidictionary.

深江 裕忠

電気通信大学大学院
情報システム学研究科

博士（工学）の学位申請論文

2014年3月

反辞書の 効率的な構築手法に 関する研究

Study on Efficient Algorithms for Construction of
an Antidictionary.

博士論文審査委員会

審査主査：	森田 啓義	教授
審査委員：	長岡 浩司	教授
審査委員：	笠井 裕之	准教授
審査委員：	古賀 久志	准教授
審査委員：	多田 好克	教授

著作権所有者：深江 裕忠, 2014

Abstract

An antidictionary is a data structure to represent minimal forbidden words which never appear on a given string but its proper substrings do appear on it.

Since Crochemore et al. originally introduced an antidictionary for a lossless data compression in 2000, it has been widely applied as a useful tool to many applications, data compression of ECG, frame-synchronization, branch prediction, and so on.

Although a naïve method for constructing an antidictionary takes time and memory complexity which is a quadratic proportional to length of a string, some algorithms with linear time and memory complexity for construction of an antidictionary have been proposed so far.

However, in spite of their linear complexity, those conventional algorithms consume a huge amount of computational time and space resources to construct an antidictionary in practical applications.

In this thesis, we discuss efficient algorithms to construct an antidictionary.

First we propose an off-line algorithm by using two array structures, that is, suffix array and L-array. A suffix array is a sorted list of the suffixes of a given string while an L-array is a list of length of the longest common prefix between adjacent pairs of suffixes on the suffix array.

Computer simulation shows that time and space complexities of the proposed algorithm are about $1/20$ and $1/2.5$ lower than a conventional linear construction algorithm, respectively.

Next, we consider a problem how to construct dynamically an antidictionary automaton which accepts minimum forbidden words of a given string.

Each time a new symbol of the string is entered into the system, we update the antidictionary automaton to fit the intermediate input data. By identifying which minimum forbidden word is deleted from or added to the current automaton, we obtain an on-line algorithm to construct an antidictionary automaton of a given string.

概要

情報技術の発展により，近年では様々な機器にセンサーが搭載され，センサーで検知したデータに従って動作する機器が多くなってきた．センサーは，温度，振動，超音波，電波など様々なものがあるが，取得するデータは数値の系列である．この数値の系列（入力系列）に対して，検出や解析を行うために，入力系列に出現する全ての部分系列を登録したデータベース（ここでは辞書と呼ぶ）がよく使われる．

本論文で研究対象にしている反辞書とは，従来の一般的な辞書とは異なり，入力系列に出現しない系列（極小禁止語）を登録する辞書である．極小禁止語とは，それ自身は入力系列に出現しないが，極小禁止語の真の部分系列は全て入力系列に出現する系列である．すなわち，入力系列に出現しない系列の中で極小性をもつ．反辞書は，全ての部分系列を登録する辞書よりも，登録する系列が少ないという性質をもつ．また，反辞書を構成するデータ構造として，極小禁止語を受理するオートマトンや，接尾辞木を拡張して極小禁止語のノードを追加するなど効率的な検索方法が提案されている．

反辞書は，Crochemoreらが2000年に提案した無歪みデータ圧縮法に初めて用いられ，符号化で使用する確率モデルの作成において，その有効性が示された．データ圧縮以外にも，反辞書は，分岐予測，同期符号，不整脈検出などへ応用されている．

しかし，反辞書を構築するには，計算量が入力系列長に比例することが知られているが，実用上は，膨大な記憶量と計算時間を必要とする．例えば，太田・森田によって提案された接尾辞木を用いた反辞書構築手法では，アルファベットサイズが多値の場合，記憶量が膨大になる．計算機実験では，アルファベットサイズが256の場合，入力系列長の約1万倍の記憶量を必要とした．

そこで，本論文では，記憶量と計算時間を改善する新しい反辞書構築手法を二つ提案する．

一つ目は，入力系列を全て入力後に反辞書を構築する静的な手法である．二つ目は，入力系列を逐次的に読み込みながら，反辞書を変化させていく動的な手法である．構築手法を考える場合，一般的には，動的な手法よりも静

的な手法の方が簡単である．そこで，最初に静的な手法について検討を行い，その成果を動的な手法へ活用する．

静的な手法では，新しい反辞書の別表現を提案する．この別表現では，極小禁止語を，先頭記号と後続する系列に分けて考えることにより，あと1記号先頭に加えることで極小禁止語になる系列が，入力系列の接尾辞配列ならびに高さ配列（本論文ではL配列と呼ぶ）によって特徴づけられること，すなわち，先頭に1記号加えれば極小禁止語になる系列は，接尾辞配列上の接尾辞の先頭部分に存在し，その長さはL配列で求めることができることを示す．ここでL配列とは，接尾辞配列上で隣接した二つの接尾辞に共通する最長接頭辞の長さの配列である．

さらに，本論文では，被覆集合を提案する．被覆集合とは，入力系列の任意の部分系列が，接尾辞配列上の接尾辞の先頭部分に存在している範囲を表現する集合である．接尾辞配列上の接尾辞は辞書式順序昇順で整列しているので，被覆集合の範囲を比較することで，入力系列の部分系列同士の辞書式順序昇順の大小関係を判別できる．すなわち，この2つの系列の先頭から順に記号を比較する必要がなく，定数時間で系列の比較が行えるようになる．

この反辞書の別表現と被覆集合を用いることで，接尾辞配列，L配列をそれぞれ3回走査するだけで，先頭に1記号加えれば極小禁止語になる系列をすべて含む高々 $2n+2$ 個（ n は入力系列長）の系列を辞書式順序昇順で求められる算法を提案する．提案方法においては，極小禁止語は，利用した被覆集合の範囲をそのまま活用して，被覆集合の範囲と系列長，そして，先頭に加える記号を要素に持つ配列で表現される．本論文では，この極小禁止語の配列表現を反辞書配列と呼んでいる．

そして，計算機実験によって提案した算法の有効性を確かめたところ，接尾辞木を用いた従来手法と比較して，計算時間が約20分の1，記憶量が約2.5分の1に改善されることが分かった．

動的な手法では，反辞書そのものではなく，反辞書オートマトン，すなわち，反辞書に含まれる極小禁止語を受理するオートマトンを，入力系列を逐次的に読み込みながら動的に構築する手法を提案する．

まず，すでに読み込まれた入力系列の末尾に新たな記号が加わることにより，反辞書がどのように更新されるかを明らかにした．申請者が修士論文で得た知見である，反辞書の更新において，ある特定の極小禁止語が一つ削除されることと，新たに加わる極小禁止語（新規極小禁止語）は，その両端のどちらか一方に削除される極小禁止語を含む，という二つの性質に加え，本論文では，新規極小禁止語の長さを削除される極小禁止語の長さで上と下か

ら評価した．

つぎに，反辞書の別表現と同様に，あと1記号加えることで新規極小禁止語になる系列を，入力系列の末尾から探索する．さらに，更新前の反辞書オートマトンを用いて，探索した系列を新規極小禁止語にするために加える記号を求める．

これらの結果を用いることによって，反辞書オートマトンの構築するために，従来の手法では，接尾辞を表す木構造やオートマトンを経由して作成する必要があったのが，その手間を省いて，入力系列から直接構築する手法を提案する．

目次

第1章	はじめに	12
1.1	研究背景	12
1.2	研究のねらい	13
第2章	定義とデータ構造	15
2.1	有限アルファベット集合と系列	15
2.2	辞書	15
2.3	接尾辞と接頭辞	16
2.4	極小禁止語と反辞書	17
2.5	辞書式順序比較	17
2.6	接尾辞配列	18
2.7	共通接頭辞の最大長	19
2.8	L 配列	19
2.9	被覆集合	20
2.10	接頭記号と接尾記号	21
第3章	接尾辞木を用いた反辞書構築	23
3.1	反辞書構築のアプローチ	23
3.2	接尾辞木	24
3.3	suffix link と MF-link	26
3.4	反辞書の構築	27
3.4.1	タイプ I の極小禁止語の追加	27
3.4.2	タイプ II の極小禁止語の追加	28
第4章	配列を用いた静的反辞書構築の提案	30
4.1	接尾辞木を用いた反辞書構築手法の記憶量の内訳	30
4.2	反辞書の別表現	30
4.3	極小禁止語の先頭文字を削除した系列	31
4.4	提案手法の概要	34

4.5	系列の表現	35
4.6	Step 3 のアルゴリズム	37
4.7	Step 4 のアルゴリズム	40
4.8	計算時間の評価	44
4.9	記憶量の評価	45
4.10	計算機実験	47
第 5 章	動的反辞書構築の提案	51
5.1	静的反辞書構築と動的反辞書構築	51
5.2	表記の拡張	51
5.3	反辞書オートマトン	52
5.4	既知の結果	54
5.5	新しく加わる極小禁止語の長さ	56
5.6	反辞書の変化と提案アルゴリズム	58
第 6 章	おわりに	61

目 次

2.1	$x = \text{MISSISSIPPI}$ のときの接尾辞配列	18
2.2	$x = \text{MISSISSIPPI}$ のときの L 配列	19
2.3	$x = \text{MISSISSIPPI}$ のときの S の被覆集合	20
2.4	$x = \text{MISSISSIPPI}$ のときの S の接尾辞配列と接頭記号	22
3.1	これまでの反辞書の構築方法	23
3.2	$x = \text{aababacbc}$ のときの接尾辞トライ (上) と接尾辞木 (下)	25
3.3	接尾辞の suffix link (実線矢印) と MF-link (点線矢印)	26
3.4	接尾辞木へのタイプ I の反辞書の追加例	28
3.5	接尾辞木へのタイプ II の反辞書の追加例	29
4.1	計算時間の比較 (実線: 提案手法, 点線: 接尾辞木を用いた手法)	48
4.2	記憶量の比較 (実線: 提案手法, 点線: 接尾辞木を用いた手法)	48
4.3	アルファベットサイズ m ごとの計算時間	50
4.4	アルファベットサイズ m ごとの記憶量	50
5.1	$x^{10} = \text{abcaababac}$ の反辞書オートマトン, 黒丸のノードが極小 禁止語を示す.	53

表 目 次

4.1	$x = \text{ababb}$ のときの \bar{t}_i と \bar{b}_i (\bar{b}_i の—は, 存在しない事を示す) . . .	33
-----	---	----

記法一覧

\mathcal{X}	有限アルファベット集合 p.15
x_i^j	系列 $x_i x_{i+1} \cdots x_j$ p.15
λ	空系列 p.15
x	長さ n の系列 $x_1 x_2 \cdots x_n$ (第4章までの記法) p.15
x^n	長さ n の系列 $x_1 x_2 \cdots x_n$ (第5章以降の記法) p.52
$\mathcal{D}(x^n)$	系列 x^n の辞書 p.15
$\mathcal{S}(x^n)$	系列 x^n の全ての接尾辞の集合 p.16
$\sigma(x^n)$	系列 x^n の先頭記号を削る接尾辞演算子 p.16
$\sigma^k(x^n)$	接尾辞演算子を系列 x^n に k 回繰り返しほどこす p.16
$\mathcal{P}(x^n)$	系列 x^n の全ての接頭辞の集合 p.16
$\pi(x^n)$	系列 x^n の末尾記号を削る接頭辞演算子 p.16
$\pi^k(x^n)$	接頭辞演算子を系列 x^n に k 回繰り返しほどこす p.16
$\mathcal{A}(x^n)$	系列 x^n の反辞書 p.17
$v \prec w$	辞書式順序昇順での大小関係 . v の方が w より小さい . . p.17
r_i	辞書式順序昇順で i 番目の接尾辞 . なお , $r_0 = \lambda$ である . p.18
S	系列 x^n の接尾辞配列 p.18
$S[i]$	系列 x^n の接尾辞配列の i ($0 \leq i \leq n$) 番目の要素 p.18
\perp	接尾辞配列の境界を示す仮想記号 p.18
$\text{lcp}(v, w)$	系列 v と w の共通接頭辞の最大の系列長 p.19

L	系列 x^n の L 配列 p.19
$L[i]$	系列 x^n の L 配列の i ($-1 \leq i \leq n$) 番目の要素 p.19
$\mathcal{C}(u)$	系列 x^n の部分系列 u の被覆集合 p.20
$\mathcal{X}^H(u)$	系列 x^n の部分系列 u の接頭記号の全ての集合 (第4章までの記法) p.21
$\mathcal{X}_n^H(u)$	系列 x^n の部分系列 u の接頭記号の全ての集合 (第5章以降の記法) p.52
$\mathcal{X}^T(u)$	系列 x^n の部分系列 u の接尾記号の全ての集合 (第4章までの記法) p.21
$\mathcal{X}_n^T(u)$	系列 x^n の部分系列 u の接尾記号の全ての集合 (第5章以降の記法) p.52
$\mathcal{N}(v)$	系列 v のローカス p.26
$L(p)$	ノード p のラベル列 p.26
$S(p)$	ノード p の suffix link 先ノード p.26
$\mathcal{H}(u)$	系列 x^n の部分系列 u に対する, $\pi(u)$ と u の全接頭記号集合の差集合 (第4章までの記法) p.30
$\mathcal{H}_n(u)$	系列 x^n の部分系列 u に対する, $\pi(u)$ と u の全接頭記号集合の差集合 p.52
$\mathcal{T}(u)$	系列 x^n の部分系列 u に対する, $\sigma(u)$ と u の全接尾記号集合の差集合 (第4章までの記法) p.31
$\mathcal{T}_n(u)$	系列 x^n の部分系列 u に対する, $\sigma(u)$ と u の全接尾記号集合の差集合 p.52
\mathcal{G}_H	系列 x^n の全ての極小禁止語の先頭文字を削除した系列の集合 . p.31
\bar{t}_i	系列 x^n の極小禁止語の先頭文字を削除した系列の候補, 長さを $L[i-1]$ で求めている p.33

\bar{b}_i	系列 x^n の極小禁止語の先頭文字を削除した系列の候補．長さを $L[i]$ で求めている．．．．．p.33
\mathcal{G}	全ての \bar{t}_i と \bar{b}_i の集合． \mathcal{G}_H を含んでいる．．．．．p.34
\mathcal{G}_T	全ての \bar{t}_i の集合．．．．．p.37
\mathcal{G}_B	全ての \bar{b}_i の集合．．．．．p.37
t_i	集合 \mathcal{G}_T の要素を辞書式順序昇順に整列したときの i ($0 \leq i \leq n$) 番目の系列．．．．．p.37
b_i	集合 \mathcal{G}_B の要素を辞書式順序昇順に整列したときの i ($0 \leq i \leq \mathcal{G}_B - 1$) 番目の系列．．．．．p.37
g_i	集合 \mathcal{G} の要素を辞書式順序昇順に整列したときの i ($0 \leq i \leq \mathcal{G} - 1$) 番目の系列．．．．．p.37

第1章 はじめに

1.1 研究背景

反辞書とは，従来の一般的な辞書とは異なり，入力系列に出現しない系列（極小禁止語）を登録する辞書である．Crochemoreらが2000年に提案した反辞書法[5]で，歪みのないデータ圧縮に有効であることが明らかにされた．データ圧縮以外にも，分岐予測[6]，同期符号[7]，不整脈検出[8]などの応用が提案されている．

反辞書を利用するには，反辞書から構築した反辞書オートマトンを用いる．例えば，反辞書法の場合，符号側では，入力系列にしたがってオートマトンの状態を遷移していき，遷移先が一意に定まるときは記号を削除することでデータを圧縮する．そして，復号側では，符号側から反辞書と生成した圧縮データが与えられる．与えられた反辞書から反辞書オートマトンを生成し，圧縮データにしたがって反辞書オートマトンを遷移することでデータを復元する．このとき，符号側で削除した記号は遷移先が一意に決まるので，データは完全に復元される．

すなわち，反辞書法は符号化側と復号側の処理がオートマトンを遷移するだけなので，計算機への負荷が少ないという利点がある．この利点は，組み込み機器などの資源が限られ，記憶量と通信路容量を十分確保できない環境下において，優れた圧縮手法となる可能性を秘めている．また，反辞書法以外についても，反辞書オートマトンを用いるので計算機への負荷が少ない．

反辞書の構築は，太田・森田によって，入力系列長に比例する，すなわち，線形の計算時間と記憶量で構築する手法[9]が提案されている．この手法では，入力記号系列の接尾辞を木構造で表現した接尾辞木を用いる．接尾辞木へMF-linkと呼ばれる特別なリンクを追加して反辞書を構築する．ただし，アルファベットが多値の場合，記憶量が膨大になるという問題がある．これは，MF-linkが使用する記憶量は線形だが，アルファベットサイズにも比例するために係数が大きくなるのが原因である．

そこで，本論文ではより少ない計算時間と記憶量で反辞書を構築すること

を検討する．

1.2 研究のねらい

本論文では，記憶量と計算時間を改善する新しい反辞書構築手法を2つ提案する．1つ目は，入力系列を全て入力後に反辞書を構築する静的な手法である．2つ目は，入力系列を逐次的に読み込みながら，反辞書を変化させていく動的な手法である．構築手法を考える場合，一般的には，動的な手法よりも静的な手法のほうが簡単である．そこで，最初に静的な手法について検討を行い，その成果を動的な手法へ活用する．

1つ目の静的手法では，接尾辞を配列で表現する接尾辞配列 [10] に注目する．ただし，単純に接尾辞木の代わりに接尾辞配列を用いると，記憶量は削減できるが，入力系列長 n に対して， n^2 に比例する計算時間を必要とする．また，接尾辞配列に高さ配列 [11],[12]（本論文では L 配列と呼ぶ）を利用した方法でも平均して $n \log n$ に比例する計算時間を必要とする問題がある．ここで L 配列とは，接尾辞配列上で隣接した二つの接尾辞に共通する最長接頭辞の長さの配列である．

そこで，本論文では計算時間を線形とするために，あと1記号を先頭に加えると反辞書の要素になる系列を接尾辞配列から探索し，反辞書を構築する手法を提案する．

そのために，極小禁止語を先頭記号と後続する系列に分けて考える，反辞書の新しい別表現を提案する．そして，あと1記号先頭に加えることで極小禁止語になる系列が，入力系列の接尾辞配列ならびに L 配列によって特徴づけられること，すなわち，先頭に1記号加えれば極小禁止語になる系列は，接尾辞配列上の接尾辞の先頭部分に存在し，その長さは L 配列で求めることができることを示す．

さらに，本論文では，被覆集合を提案する．被覆集合とは，入力系列の任意の部分系列が，接尾辞配列上の接尾辞の先頭部分に存在している範囲を表現する集合である．接尾辞配列上の接尾辞は辞書式順序昇順で整列しているので，被覆集合の範囲を比較することで，入力系列の部分系列同士の辞書式順序昇順の大小関係を判別できる．すなわち，この2つの系列の先頭から順に記号を比較する必要がなく，定数時間で系列の比較が行えるようになる．

この反辞書の別表現と被覆集合を用いることで，接尾辞配列， L 配列をそれぞれ3回走査するだけで，先頭に1記号加えれば極小禁止語になる系列をすべて含む高々 $2n + 2$ 個の系列を辞書式順序昇順で求められる算法を提案す

る．提案方法においては，極小禁止語は，利用した被覆集合の範囲をそのまま活用して，被覆集合の範囲と系列長，そして，先頭に加える記号を要素に持つ配列で表現される．本論文では，この極小禁止語の配列表現を反辞書配列と呼んでいる．

また，計算機実験を行い，接尾辞木を用いた従来手法と比較して，計算時間が約 20 分の 1，記憶量が約 2.5 分の 1 に改善されることを示す．

2 つ目の動的な手法では，反辞書そのものではなく，反辞書オートマトン，すなわち，反辞書に含まれる極小禁止語を受理するオートマトンを，入力系列を逐次的に読み込みながら動的に構築する手法を提案する．

そのために，入力系列の末尾に新たな記号が加わることで反辞書がどのように更新するのかを検討する．

これまでの研究で，反辞書が変化するときには，削除される極小禁止語と新たに加える極小禁止語（新規極小禁止語）は明らかになっている．しかし，効率的に求めることができない問題があった．本論文では，さらに検討を進め，新たに加わる極小禁止語の長さを評価する．

つぎに，反辞書の別表現と同様に，あと 1 記号加えることで新規極小禁止語になる系列を，入力系列の末尾から探索する．さらに，更新前の反辞書オートマトンを用いて，探索した系列を新規極小禁止語にするために加える記号を求める．

これらの結果を用いることによって，反辞書オートマトンの構築するために，従来の手法では，接尾辞を表す木構造やオートマトンを経由して作成する必要があったのが，その手間を省いて，入力系列から直接構築する手法を提案する．

第2章 定義とデータ構造

2.1 有限アルファベット集合と系列

有限アルファベット集合を

$$\mathcal{X} = \{\xi_1, \xi_2, \dots, \xi_m\} \ (\xi_1 < \xi_2 < \dots < \xi_m, m \geq 2)$$

とおき, \mathcal{X} 上の系列を,

$$x_i^j = x_i x_{i+1} \cdots x_j \ (x_k \in \mathcal{X}, 1 \leq i \leq k \leq j)$$

とする.

ここで, 集合のサイズと系列長は $|\cdot|$ で表記する. 例えば, $|\mathcal{X}| = m, |x_i^j| = j - i + 1$ である. なお, 系列長 0 を空系列とし, λ と表記する. 便宜上, $i > j$ のとき $x_i^j = \lambda$ とする. また, 系列長 $n (\geq 1)$ の系列 x_1^n を x と表記する.

2.2 辞書

辞書 (dictionary) $\mathcal{D}(x)$ とは, x 上に出現する部分系列の全ての集合であり,

$$\mathcal{D}(x) = \{x_i^j | 1 \leq i \leq j \leq n\} \cup \{\lambda\}$$

とする. なお, $\mathcal{D}(\lambda) = \{\lambda\}$ であることを注記する.

例 1 系列 $x = abbca$ の辞書は,

$$\mathcal{D}(x) = \{a, b, c, ab, bb, bc, ca, abb, bbc, bca, abbc, bbca, abbca, \lambda\}$$

となる.

2.3 接尾辞と接頭辞

系列 $x = x^n$ に対して x の接尾辞 (suffix) とは, x_i^n ($1 \leq i \leq n$) に空系列 λ を加えたものである. ここで, x の全ての接尾辞からなる集合 $S(x)$ は,

$$S(x) = \{x_i^n | 1 \leq i \leq n\} \cup \{\lambda\}$$

と表される. さらに, 接尾辞演算子 σ を $\sigma(x) = x_2^n$ と定義する. x に σ を k (≥ 1) 回繰り返し施すと,

$$\sigma^k(x) = x_{1+k}^n$$

となる. すなわち, $\sigma^k(x)$ は, x の先頭から k 記号分削除した系列である.

また, x_1^i ($1 \leq i \leq n$) と空系列 λ を x の接頭辞 (prefix) と呼ぶ. ここでは, x の全ての接頭辞からなる集合 $P(x)$ は,

$$P(x) = \{x_1^i | 1 \leq i \leq n\} \cup \{\lambda\}$$

と表される. さらに, 接頭辞演算子 π を $\pi(x) = x_1^{n-1}$ と定義する. x に π を k (≥ 1) 回繰り返し施すと,

$$\pi^k(x) = x_1^{n-k}$$

となる. すなわち, $\pi^k(x)$ は, x の末尾から k 記号分削除した系列である.

なお, $\sigma(\lambda) = \lambda, \pi(\lambda) = \lambda$ であることを注記する.

例 2 系列 $x = abbca$ の全ての接尾辞の集合 $S(x)$ と, 全ての接頭辞の集合 $P(x)$ は,

$$S(x) = \{abbca, bbca, bca, ca, a, \lambda\},$$

$$P(x) = \{abbca, abbc, abb, ab, a, \lambda\}$$

となる.

また,

$$\sigma(x) = bbca, \quad \sigma^2(x) = bca,$$

$$\pi(x) = abbc, \quad \pi^2(x) = abb$$

となる.

2.4 極小禁止語と反辞書

系列 x 上に出現しない系列を禁止語とよぶ．そして，系列 v が次の (2.1), (2.2), (2.3) のすべてを満すとき， v は x の極小禁止語 (Minimal Forbidden Word) である．

$$v \notin \mathcal{D}(x), \quad (2.1)$$

$$\pi(v) \in \mathcal{D}(x), \quad (2.2)$$

$$\sigma(v) \in \mathcal{D}(x). \quad (2.3)$$

したがって，極小禁止語は x 上に出現しないが，極小禁止語の端を削除した系列は x 上に出現している．

反辞書 (antidictionary) $\mathcal{A}(x)$ とは，系列 x に対する極小禁止語のすべての集合である．すなわち，

$$\mathcal{A}(x) = \{v \notin \mathcal{D}(x) \mid \pi(v) \in \mathcal{D}(x), \sigma(v) \in \mathcal{D}(x)\}$$

となる．

例 3 有限アルファベット集合 $\mathcal{X} = \{a, b, c\}$ ，系列 $x = abbca$ の反辞書は，

$$\mathcal{A}(x) = \{aa, ac, ba, cb, cc, abc, bbb, cab\}$$

となる．

2.5 辞書式順序比較

系列 v, w の辞書式順序昇順での大小関係を，比較記号 \prec で示すことにする． $v \prec w$ は，

$$(v_1 < w_1) \vee (v_1 = w_1) \wedge (\sigma(v) \prec \sigma(w))$$

と定義する．便宜上， $\lambda < \xi_1$ とする．

例 4 次の 5 つの系列 $bb, aca, c, abba, abb$ の辞書式順序昇順での大小関係は，

$$abb \prec abba \prec aca \prec bb \prec c$$

となる．

2.6 接尾辞配列

系列 x の全ての接尾辞を辞書式順序昇順で並べたとき, i ($0 \leq i \leq n$) 番目の接尾辞をランク i の接尾辞と呼び, r_i と表記する. すなわち, $r_0 = \lambda$ となり, $i > 0$ のときは,

$$r_i = \min\{v \in \mathcal{S}(x) \mid r_{i-1} \prec v\} \quad (2.4)$$

となる.

系列 x の接尾辞配列 (suffix array) S とは, S の i 番目の要素が

$$S[i] = n - |r_i| + 1$$

与えられる配列である. この $S[i]$ は接尾辞 r_i の x 上におけるインデックスを表す. よって, x の左端から $S[i] - 1$ 個の記号を削ると, r_i になるので,

$$r_i = \sigma^j(x) \quad (j = S[i] - 1)$$

が成り立つ.

なお, 便宜上, 本論では境界を示すために, \mathcal{A} に属さない記号 \perp を用いて $r_{-1} = r_{n+1} = \perp$ とおく.

例 5 図 2.1 に, $x = \text{MISSISSIPPI}$ (系列長 11) の接尾辞配列の概要を示す.

最初に, x の全ての接尾辞を辞書式順序昇順で整列する. 整列後の並びで i ($0 \leq i \leq 11$) 番目の接尾辞が r_i である. また, r_i の x 上の開始インデックスが $S[i]$ の値となる.

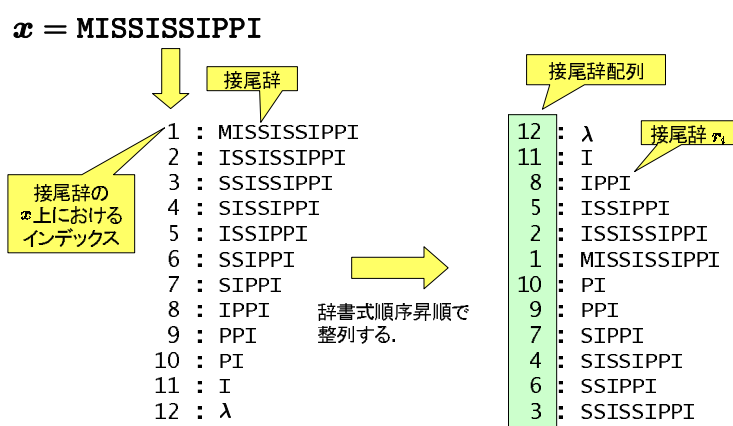


図 2.1: $x = \text{MISSISSIPPI}$ のときの接尾辞配列

2.7 共通接頭辞の最大長

2つの任意の系列 $v, w \in \mathcal{X}^*$ には, 共通接頭辞 $u \in \mathcal{P}(v) \cap \mathcal{P}(w)$ がある. この共通接頭辞の最大の系列長を

$$\text{lcp}(v, w) = \max\{|u| \mid u \in \mathcal{P}(v) \cap \mathcal{P}(w)\}$$

とする. 記号 \perp は \mathcal{X} に属さないので, v との共通接頭辞は存在しない. そこで便宜上, $\text{lcp}(\perp, v) = \text{lcp}(v, \perp) = -1$ とする.

2.8 L 配列

L 配列は高さ配列 [11],[12] と呼ばれ, 2つの接尾辞 r_i, r_{i+1} ($-1 \leq i \leq n$) の最長共通接頭辞の系列長を配列で表現したものである.

すなわち, x の L 配列 L とは, L の i 番目の要素が

$$L[i] = \text{lcp}(r_i, r_{i+1}) \quad (2.5)$$

で与えられる配列である. なお, $L[-1] = L[n] = -1$ で, L の長さが $n+2$ であることを注記する.

例 6 図 2.2 に, $x = \text{MISSISSIPPI}$ (系列長 11) の L 配列の概要を示す.

r_3 と r_4 の最長共通接頭辞は ISSI で, 系列長は 4 である. よって, $L[3] = 4$ となる. 便宜上, $L[-1] = L[n] = -1$ は境界を示すために存在する.

$x = \text{MISSISSIPPI}$

$L[i] = \text{lcp}(r_i, r_{i+1})$

接尾辞配列上で
下の接尾辞との
共通接頭辞の最大長

i	$S[i]$	r_i	$L[i]$
-1	\perp		-1
0	12		0
1	11	I	1
2	8	IPPI	1
3	5	ISSIPPI	4
4	2	ISSISSIPPI	0
5	1	MISSISSIPPI	0
6	10	PI	1
7	9	PPI	0
8	7	SIPPI	2
9	4	SISSIPPI	1
10	6	SSIPPI	3
11	3	SSISSIPPI	-1
\perp	\perp		\perp

図 2.2: $x = \text{MISSISSIPPI}$ のときの L 配列

2.9 被覆集合

系列 x 上の任意の部分系列 $u \in \mathcal{D}(x)$ に対して, u を接頭辞にもつ接尾辞が必ず存在する. そこで,

$$\mathcal{C}(u) = \{i | u \in \mathcal{P}(r_i), 0 \leq i \leq n\} \quad (2.6)$$

を定義し, 被覆集合(cover set)と呼ぶことにする. 接尾辞が整列しているから, $\mathcal{C}(u)$ は連続した整数からなる集合である. すなわち,

$$\mathcal{C}(u).min = \min\{i | i \in \mathcal{C}(u)\}, \quad (2.7)$$

$$\mathcal{C}(u).max = \max\{i | i \in \mathcal{C}(u)\} \quad (2.8)$$

とすると,

$$\mathcal{C}(u) = \{i | \mathcal{C}(u).min \leq i \leq \mathcal{C}(u).max\} \quad (2.9)$$

と表せる. なお, 被覆集合は, 文献[12]の rank interval を整数の集合で表現したものである. ここで, $|\mathcal{C}(u)|$ は u が x 上に出現する数に等しく, また, $\mathcal{C}(\lambda) = \{0, 1, 2, \dots, n\}$ であることを注記する.

例 7 図 2.3 に, $x = \text{MISSISSIPPI}$ (系列長 11) の被覆集合の概要を示す.

部分系列 s の被覆集合を求めるとき, 接尾辞の先頭部分に s があるものを探す. この場合, 接尾辞 r_8, r_9, r_{10}, r_{11} が該当する. したがって, $\mathcal{C}(s) = \{8, 9, 10, 11\}$ となる.

任意の部分系列を
接頭辞にもつ接尾辞が
必ず存在する

Sの被覆集合

$\mathcal{C}(s) = \{8, 9, 10, 11\}$

i	$ S[i] $	r_i
0	12	λ
1	11	I
2	8	I P P I
3	5	I S S I P P I
4	2	I S S I S S I P P I
5	1	M I S S I S S I P P I
6	10	P I
7	9	P P I
8	7	S I P P I
9	4	S I S S I P P I
10	6	S S I P P I
11	3	S S I S S I P P I

図 2.3: $x = \text{MISSISSIPPI}$ のときの s の被覆集合

2.10 接頭記号と接尾記号

系列 $u \in \mathcal{D}(x)$ と任意の記号 $\alpha \in \mathcal{X}$ に対して, x 上に αu が出現するなら, α を u の接頭記号(head symbol) と呼ぶことにする. さらに, u の接頭記号の全ての集合を, u の全接頭記号集合と呼ぶことにし, $\mathcal{X}^H(u)$ と表記する. すなわち,

$$\mathcal{X}^H(u) = \{\alpha \in \mathcal{X} | \alpha u \in \mathcal{D}(x)\} \quad (2.10)$$

である.

さらに, $\mathcal{X}^H(u)$ は, 接尾辞配列 S と被覆集合 $\mathcal{C}(u)$ を用いて

$$\mathcal{X}^H(u) = \{x_i | i = S[k] - 1, k \in \mathcal{C}(u), i \geq 1\} \quad (2.11)$$

と表すこともできる.

また, x 上に $u\beta$ が出現するなら, β を u の接尾記号(tail symbol) と呼ぶことにする. さらに, u の接尾記号の全ての集合を, u の全接尾記号集合と呼ぶことにし, $\mathcal{X}^T(u)$ と表記する. すなわち,

$$\mathcal{X}^T(u) = \{\beta \in \mathcal{X} | u\beta \in \mathcal{D}(x)\} \quad (2.12)$$

である.

なお, $\mathcal{X}^H(\lambda)$ と $\mathcal{X}^T(\lambda)$ は x 上に出現する全ての記号であることを注記する. また, u が x 上に1箇所のみ出現し, $u \in \mathcal{P}(x)$ のときは, $\mathcal{X}^H(u) = \phi$ である. 同じく, u が x 上に1箇所のみ出現し, $u \in \mathcal{S}(x)$ のときは, $\mathcal{X}^T(u) = \phi$ である.

例 8 図 2.4 に, $x = \text{MISSISSIPPI}$ の接尾辞配列と接頭記号に関する概要を示す.

接尾辞 $r_3 = \text{ISSIPPI}$ の x 上における開始インデックスは $S[3] = 5$ である. すなわち, 記号 $x_{S[3]-1} = s$ は, r_3 の接頭記号である. さらに, r_3 の接頭辞についても, 開始インデックスは変わらないので, 記号 $x_{S[3]-1}$ は, r_3 の接頭辞の接頭記号である.

全接頭記号集合は, 被覆集合を使うことで接尾辞配列から求めることができる. 部分系列 s の被覆集合は, $\mathcal{C}(s) = \{8, 9, 10, 11\}$ なので, r_8, r_9, r_{10}, r_{11} の接頭記号を集めたものが全接頭記号集合 $\mathcal{X}^H(s)$ になる.

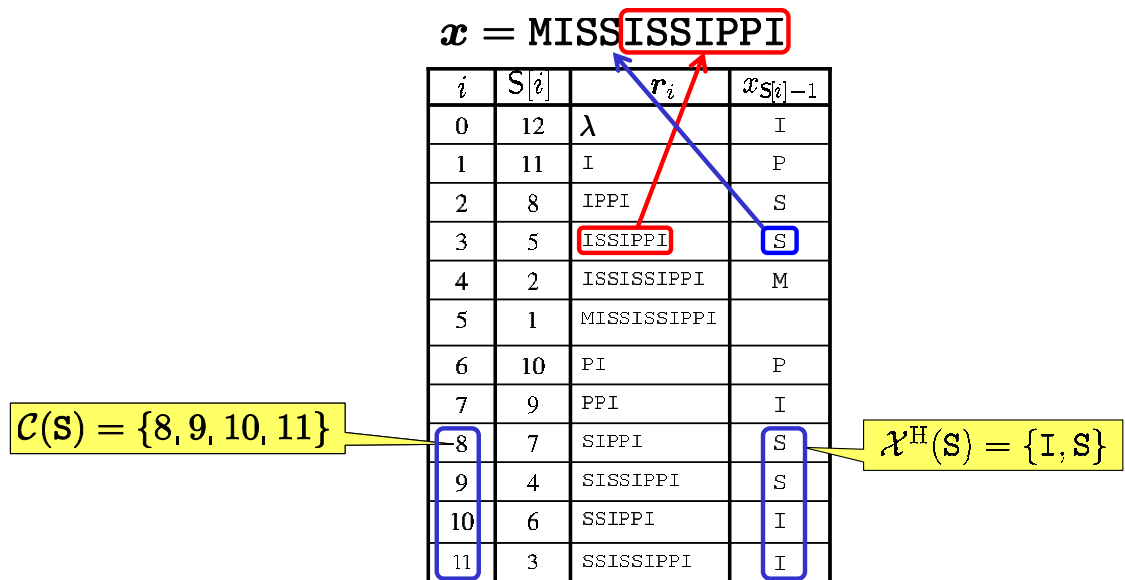


図 2.4: $x = \text{MISSISSIPPI}$ のときの S の接尾辞配列と接頭記号

第3章 接尾辞木を用いた反辞書構築

3.1 反辞書構築のアプローチ

反辞書構築の従来手法では，入力系列から直接反辞書を構築するのではなく，一度接尾辞を用いたデータ構造を構築し，そこから反辞書を構築している．また，入力系列を全て入力してから反辞書を構築する静的構築と，入力系列を1記号ずつ読み込みながら逐次的に反辞書を変化させて構築する動的構築の2種類がある（図3.1）．

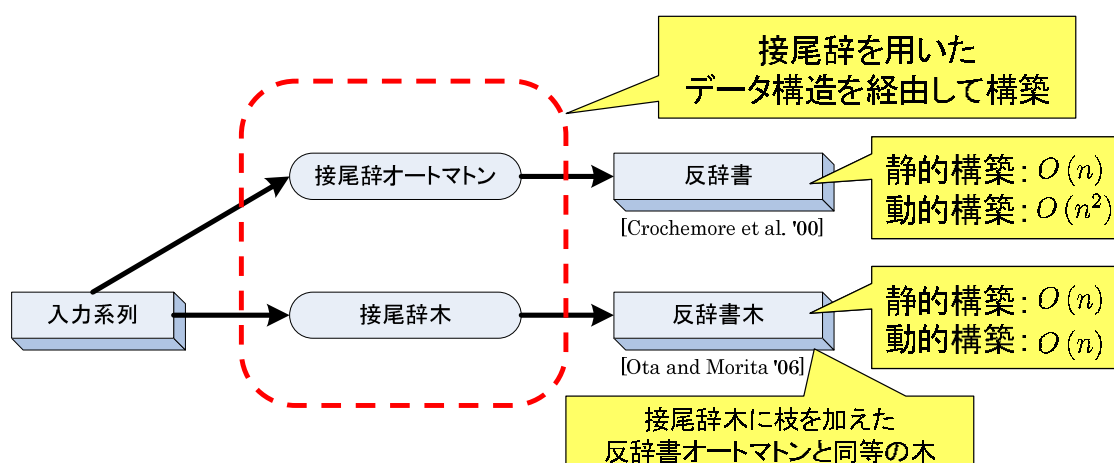


図 3.1: これまでの反辞書の構築方法

接尾辞を用いたデータ構造の多くは入力系列をすべて読み込み終わってから構築するアルゴリズムが多いため，反辞書を静的に構築する場合には，これらのデータ構造を活用できる利点がある．ただし，データ圧縮などに応用するときには，構築した反辞書を復号側に送るなどのオーバーヘッドが生じる．

一方，反辞書を動的に構築する場合には，入力記号に合わせながら反辞書を変化させていくので，オーバーヘッドを無くすることができる．また，1記号読み込んだ後の反辞書の変化が少ない場合，第4章の提案手法よりも効率が高くなる可能性がある．

第4章と第5章では、静的構築と動的構築の両方で入力系列長に対して線形の計算量を達成している、太田・森田によって提案された接尾辞木を用いた反辞書構築手法 [9] を参考にして、新しい構築手法を提案した。

そこで、この章では接尾辞木を用いた反辞書構築手法について述べる。

3.2 接尾辞木

系列の集合を木構造で表現したものとして、トライ (trie) [13] がある。トライの枝にはラベルが付けられ、記号が1つ割り振られている。そして、与えられた系列 x の全ての接尾辞の集合 $S(n)$ をトライで表現した木を接尾辞トライと呼ぶ。枝のラベルには、接尾辞上の記号が1つ割り振られている。

一方、接尾辞木 (Suffix tree) [14] も同様に、 $S(n)$ を木構造で表現したものがある。接尾辞トライと異なるのは、ラベルには x の部分系列を割り振っていることである。ここで、接尾辞トライのノードのうち、子ノードを1つだけ持つノードを陰ノード、子ノードを2つ以上持つノードを陽ノードと呼ぶことにする。任意の陰ノードを p 、 p の親ノードを q 、 p の子ノードを r とし、 p と q を接続する枝のラベルを α 、 p と r を接続する枝のラベルを β とする。このとき、系列 $\alpha\beta$ は x の部分系列である。したがって、枝に x の部分系列を割り振れるのならば、 p をなくして、 q と r を1本の枝 (ラベルは $\alpha\beta$) にまとめることができる。このようにして、接尾辞トライの陰ノードを1本の枝にまとめ、接尾辞トライをコンパクトに表現したものが接尾辞木である。

ところで、接尾辞木の枝のラベルの長さが2以上の時、枝を記号単位で分割した区切り位置に仮想のノードがあると考えれば接尾辞トライと等価な木になる。このとき、この仮想ノードは接尾辞トライの陰ノードに対応し、接尾辞木の間接ノードは接尾辞トライの陽ノードに対応する。本論文では、接尾辞木において、上記の仮想ノードを扱うことがあるので、この仮想ノードを陰ノードと呼ぶことにする。

また、陽ノードから子ノードに向かって、枝が伸びている。これらの枝のラベルの先頭記号は、各々異なる記号である。そこで、陽ノードから子ノードへ記号 α を先頭に持つラベルの枝が伸びているとき、 α が発芽すると呼ぶことにする。

図3.2に $x = \text{aababacbc}b$ の接尾辞トライと接尾辞木を示す。

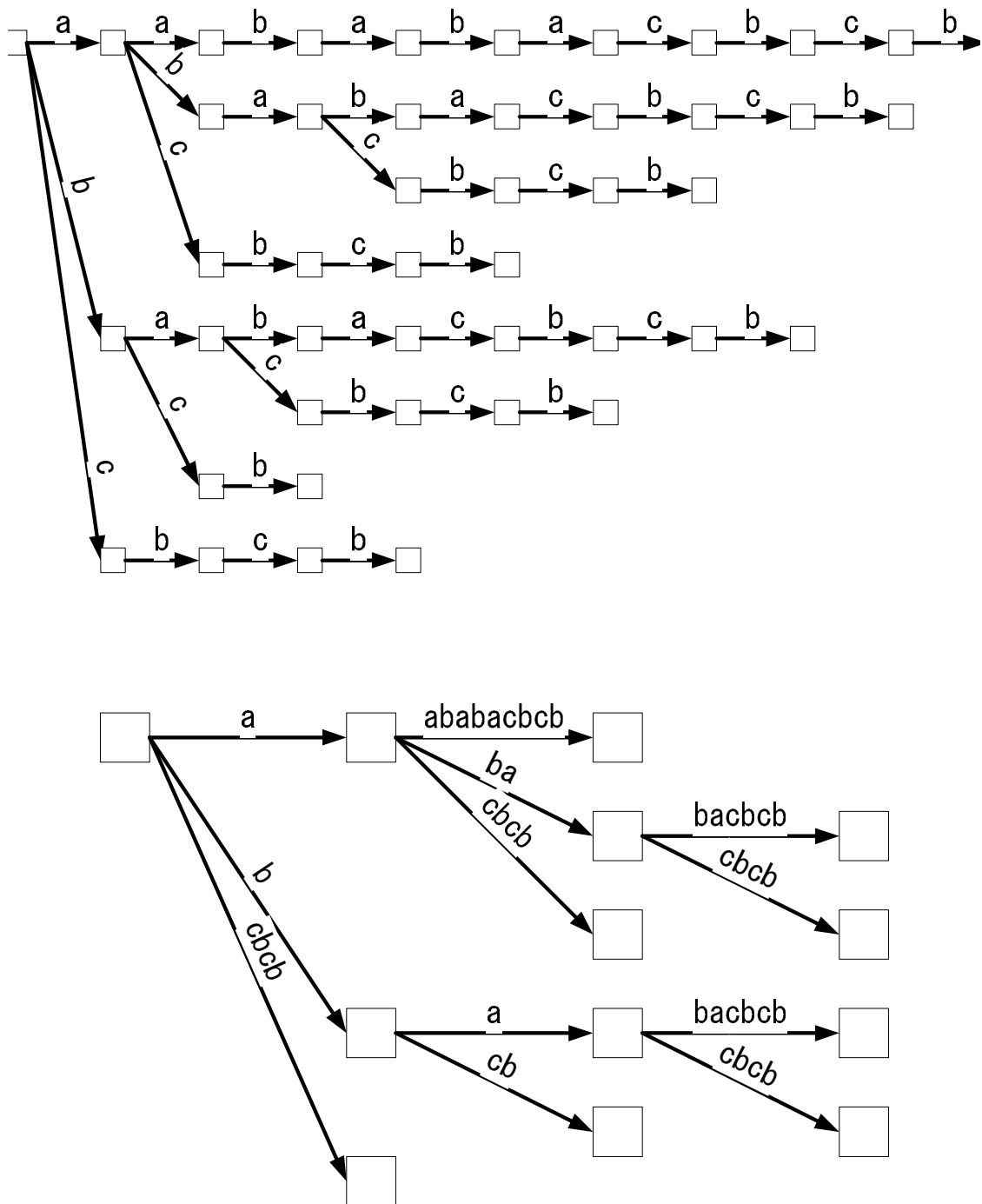


図 3.2: $x = \text{aababacbcbb}$ のときの接尾辞トライ (上) と接尾辞木 (下)

また，任意の系列 $v \in \mathcal{X}^*$ に対して，接尾辞木 \mathbb{T} の根から葉ノードに向かって v と一致するようにラベルをたどることで到達したノードをローカスと呼び， $\mathcal{N}(v)$ と表記する．系列 v によっては， $\mathcal{N}(v)$ は陰ノードだったり，存在しないこともありえる． $\mathcal{N}(v)$ が存在するときは， $v \in \mathcal{D}(x)$ である．

3.3 suffix link と MF-link

また, 太田らは, 接尾辞木を用いて反辞書を作成するために, MF-linkを定義した[9]. これは, ノード p の逆 suffix link である. すなわち, ノード m のラベル列 $L(m)$ が $\sigma(L(m)) = L(p)$ を満たすとき, p から m へ MF-link が張られる. なお, m は陽ノードだけでなく陰ノードのものもあることを注記する.

26

3.4 反辞書の構築

接尾辞木を用いた反辞書構築手法では、極小禁止語をラベル列に持つMFWノードを接尾辞木に追加する。そして、MFWノードを追加する位置によって、極小禁止語を2種類に分けている。

- 葉ノードにMFWノードを追加する極小禁止語（タイプI）。
- 葉ノード以外にMFWノードを追加する極小禁止語（タイプII）。

Ukkonenの接尾辞木の構築手法では、1記号ずつ読み込みながら接尾辞木を変化させている。1記号の読み込みにつき、葉ノードと陽ノードが高々1つ増える。したがって、接尾辞木の変化に追従しながら、上記の2種類の極小禁止語を求めていくことで、反辞書を構築する。

3.4.1 タイプIの極小禁止語の追加

タイプIの極小禁止語の場合、接尾辞木の葉ノードにMFWノードを追加する。すなわち、葉ノードを q とおいたとき、極小禁止語は $L(q)\beta$ ($\beta \in \mathcal{X}$)となる。このとき、 $L(q) \in \mathcal{D}(x)$ なので、極小禁止語の3つの条件のうち、(2.2)を満たしている。さらに、 q は葉ノードなので、どんな β でも、 $L(q)\beta \notin \mathcal{D}(x)$ となり、極小禁止語の3条件のうち、(2.1)を満たす。よって、(2.3)を満たす β を求めればよい。

ここで、 β は q のsuffix linkを用いて求める。 q のsuffix link先 $S(q)$ のラベル列は $\sigma(L(q))$ である。すなわち、 $S(q)$ から発芽する記号を β とすれば、 $\sigma(L(q))\beta \in \mathcal{D}(x)$ となり、(2.3)を満たす。

こうして、ノード q から、求めた β を枝のラベルにしてMFWノードを伸ばして追加する。

なお、ラベル列が最短でない葉ノードのsuffix link先は、別の葉ノードである。したがって、suffix link先で発芽する記号はない。言い換えると、タイプIの極小禁止語のMFWノードは、ラベル列が最短の葉ノードだけに追加する。

例9 図3.4に、 $x = \text{abbca}$ のタイプIのMFWノードを追加するときの様子を示す。図の白丸のノードが接尾辞木の陽ノード、四角のノードが接尾辞木の葉ノードである。三角のノードが追加するMFWノードである。

ラベル列が最短の葉ノードを q とすると、 $L(q) = \text{ca}$ である。すなわち、タイプIの極小禁止語は $\text{ca}\beta$ ($\beta \in \mathcal{X}$)とおける。

ここで, q の suffix link 先 $S(q)$ のラベル列は, a である. また, $S(q)$ からは, b が発芽しているので, $ab \in \mathcal{D}(x)$ である. よって, β が b のときは, cab が極小禁止語の 3 条件を満たすので, 葉ノード q から MFW ノードを伸ばしている.

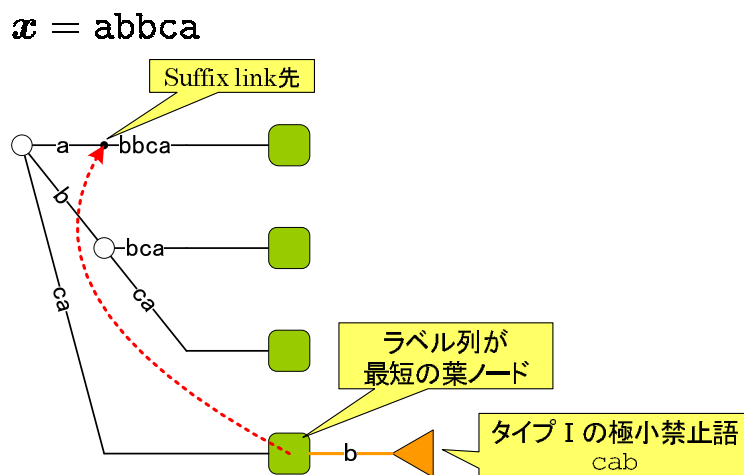


図 3.4: 接尾辞木へのタイプ I の反辞書の追加例

3.4.2 タイプ II の極小禁止語の追加

タイプ II の極小禁止語の場合, 接尾辞木の葉ノード以外のノード, つまり陽ノードや陰ノードに MFW ノードを追加する. 追加する位置は, 陽ノード p から張られている MF-link 先 m である. すなわち, 極小禁止語は $L(m)\beta$ となる. このとき, $L(m) \in \mathcal{D}(x)$, なので, 極小禁止語の 3 つの条件のうち, (2.2) を満たしている. よって, (2.1) と (2.3) を満たす β を求めればよい.

ここで, p と m から発芽する記号に注目する. p は陽ノードなので, 複数の記号が発芽している. そこで, p で発芽し, m で発芽しない記号を β とすれば, $L(m)\beta \notin \mathcal{D}(x)$, $\sigma(L(m)\beta) = L(p)\beta \in \mathcal{D}(x)$ となり, (2.1) と (2.3) を満たす.

こうして, ノード m から, 求めた β を枝のラベルにして MFW ノードを伸ばして追加する.

例 10 図 3.5 に, $x = abbca$ のタイプ II の MFW ノードを追加するときの様子を示す.

ラベル列が b の陽ノードを p , ノード p の MF-link 先のうち, ラベル列が ab のノードを m とする. このとき, タイプ II の極小禁止語は $ab\beta$ ($\beta \in \mathcal{X}$) とおける.

ノード p から発芽している記号は b と c , ノード m から発芽している記号は b である . よって , β が c のときは , abc が極小禁止語の 3 条件を満たすので , ノード m から MFW ノードを伸ばしている .

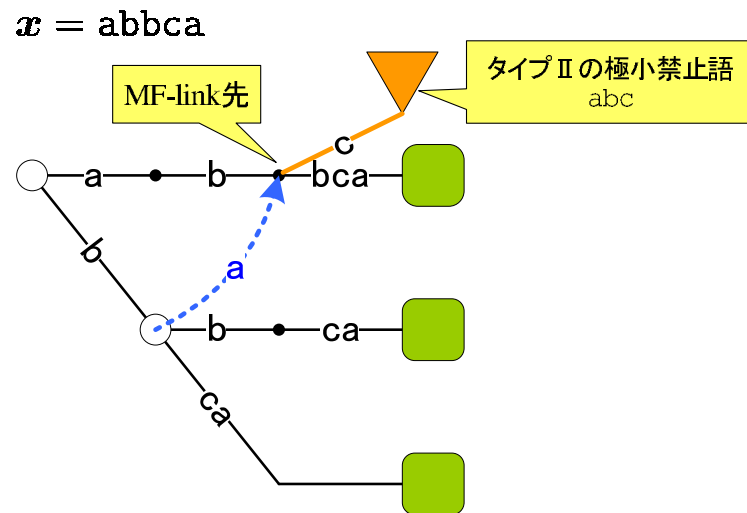


図 3.5: 接尾辞木へのタイプ II の反辞書の追加例

第4章 配列を用いた静的反辞書構築の提案

4.1 接尾辞木を用いた反辞書構築手法の記憶量の内訳

接尾辞木を用いた反辞書構築手法は、計算量が線形とはいえ、アルファベットが多値の場合には係数が大きくなるという問題がある。この問題を解決するために、使用している記憶量の内訳を調べると、記憶量はMF-linkを保持するためにほとんど占められていることがわかった。これは、MF-linkが使用する記憶量がアルファベットサイズにも比例するのが原因である。

この問題を解決するためには、MF-linkを保持しないで反辞書を構築する新たな手法が必要となる。

そこで本研究では、接尾辞木よりも少ない記憶量で接尾辞を表現できる接尾辞配列に注目した。この章では、接尾辞配列を用いて、MF-linkを保持しない新たな反辞書構築手法を提案する。

4.2 反辞書の別表現

接尾辞配列を用いて極小禁止語を求めるために、極小禁止語を先頭記号 α と系列 u の2つに分割することを提案する [1], [2]。

そのために、(2.10) で定義した全接頭記号集合 $\mathcal{X}^H(\cdot)$ を用いて、 $\mathcal{X}^H(\pi(u))$ と $\mathcal{X}^H(u)$ との差集合を $\mathcal{H}(u)$ と表記し、

$$\mathcal{H}(u) = \begin{cases} \mathcal{X} \setminus \mathcal{X}^H(\lambda) & (u = \lambda) \\ \mathcal{X}^H(\pi(u)) \setminus \mathcal{X}^H(u) & (u \neq \lambda) \end{cases} \quad (4.1)$$

と定義する。

そして、 $\mathcal{H}(u)$ を用いて反辞書 $\mathcal{A}(x)$ を表せることを、次の定理 1 で示す。

定理 1 任意の $x \in \mathcal{X}^*$ に対して、

$$\mathcal{A}(x) = \{\alpha u \mid \alpha \in \mathcal{H}(u), u \in \mathcal{D}(x)\}.$$

証明 最初に, $\alpha u \in \mathcal{A}(x)$ ($\alpha \in \mathcal{X}, u \in \mathcal{D}(x)$) とおいたとき, $\alpha \in \mathcal{H}(u)$ であることを, $u = \lambda$ と $u \neq \lambda$ の場合に分けて証明する.

まず, $u = \lambda$ の場合, $\alpha \in \mathcal{A}(x)$ となるので, α は x 上に出現していない記号である. ここで, $\mathcal{X}^H(\lambda)$ は x 上に出現する全ての記号の集合なので, $\alpha \in \mathcal{X} \setminus \mathcal{X}^H(\lambda) = \mathcal{H}(u)$ である.

また, $u \neq \lambda$ の場合, $\alpha u \in \mathcal{A}(x)$ から, $\alpha u \notin \mathcal{D}(x), \alpha \pi(u) \in \mathcal{D}(x)$ となる. (2.10) から, $\alpha \notin \mathcal{X}^H(u), \alpha \in \mathcal{X}^H(\pi(u))$ となる. よって, $\alpha \in \mathcal{X}^H(\pi(u)) \setminus \mathcal{X}^H(u) = \mathcal{H}(u)$ である.

次に, $u \in \mathcal{D}(x), \alpha \in \mathcal{H}(u)$ ならば, $\alpha u \in \mathcal{A}(x)$ であることを, $u = \lambda$ と $u \neq \lambda$ の場合に分けて証明する.

まず, $u = \lambda$ の場合, $\mathcal{X}^H(\lambda)$ は x 上に出現する全ての記号の集合である. よって, (4.1) より, α は未出現記号である. このとき, αu は, (2.1), (2.2), (2.3) を満す. よって, $\alpha u \in \mathcal{A}(x)$ である.

また, $u \neq \lambda$ の場合, (4.1) から, $\alpha \in \mathcal{X}^H(\pi(u)), \alpha \notin \mathcal{X}^H(u)$ である. このとき (2.10) から, $\alpha \pi(u) \in \mathcal{D}(x), \alpha u \notin \mathcal{D}(x)$ となり, (2.1) と (2.2) を満す. また, $u \in \mathcal{D}(x)$ なので, (2.3) を満す. よって, $\alpha u \in \mathcal{A}(x)$ である.

以上より, 題意が成り立つ.

□

ところで, (2.12) で定義した全接尾記号集合 $\mathcal{X}^T(\cdot)$ を用いて,

$$\mathcal{X}^T(u) = \{\beta \in \mathcal{X} \mid u\beta \in \mathcal{D}(x)\}, \quad (4.2)$$

$$\mathcal{T}(u) = \begin{cases} \mathcal{X} \setminus \mathcal{X}^T(\lambda) & (u = \lambda) \\ \mathcal{X}^T(\sigma(u)) \setminus \mathcal{X}^T(u) & (u \neq \lambda) \end{cases} \quad (4.3)$$

と定義すると, 定理 1 と同様に次の系 1 が成り立つ.

系 1 任意の $x \in \mathcal{X}^*$ に対して,

$$\mathcal{A}(x) = \{u\beta \mid \beta \in \mathcal{T}(u), u \in \mathcal{D}(x)\}.$$

なお, 系 1 は, この章では使わないが, 第 5 章で用いる.

4.3 極小禁止語の先頭文字を削除した系列

任意の部分系列 $u \in \mathcal{D}(x)$ が $\mathcal{H}(u) \neq \phi$ のとき,

$$\mathcal{G}_H = \{u \mid u \in \mathcal{D}(x), \mathcal{H}(u) \neq \phi\} \quad (4.4)$$

と定義する．そして，

$$\sigma(\mathcal{A}(x)) = \{\sigma(w) | w \in \mathcal{A}(x)\}$$

とおくと，定理 1 より，

$$\mathcal{G}_H = \sigma(\mathcal{A}(x))$$

が成立する．

この \mathcal{G}_H を求めることができれば，定理 1 から，要素 $u \in \mathcal{G}_H$ に対して $\mathcal{H}(u)$ を求めることで反辞書を構築できる．しかし， \mathcal{G}_H を求めることは容易ではない．そこで，容易に求めれて， \mathcal{G}_H を含む集合を考えることにする．

まずは，(2.4) で定義したランク i の接尾辞 r_i ($0 \leq i \leq n$) ごとに， r_i の全ての接頭辞の集合 $\mathcal{P}(r_i)$ を 2 つに分割した集合 $\mathcal{F}(i), \mathcal{E}(i)$ を定義する．

$$\mathcal{F}(i) = \{u \in \mathcal{P}(r_i) | |u| \neq \text{lcp}(r_i, r_k) + 1, -1 \leq \forall k \leq n+1\}, \quad (4.5)$$

$$\mathcal{E}(i) = \{u \in \mathcal{P}(r_i) | |u| = \text{lcp}(r_i, r_k) + 1, -1 \leq \exists k \leq n+1\}. \quad (4.6)$$

明らかに， $\mathcal{P}(r_j) = \mathcal{F}(j) \cup \mathcal{E}(j), \mathcal{F}(j) \cap \mathcal{E}(j) = \phi$ となる．

次の補題 1 で， $\mathcal{F}(i)$ と \mathcal{G}_H の関係を示す．

補題 1 任意の $v \in \mathcal{X}^*$ とある i ($0 \leq i \leq n$) に対して， $v \in \mathcal{F}(i)$ ならば，

$$\mathcal{H}(v) = \phi.$$

証明 空系列 λ は， $\lambda \in \mathcal{E}(i)$ である．よって， $v \neq \lambda$ ．そこで， $v = u\beta$ ($u = \pi(v), \beta \in \mathcal{X}$) とおく．

ここで，背理法を使って証明する．いま， $\mathcal{H}(v) \neq \phi$ と仮定する．このとき，ある $\alpha \in \mathcal{H}(v)$ が存在し， $\alpha \in \mathcal{X}^H(u) \setminus \mathcal{X}^H(u\beta)$ とおける．すなわち， $\alpha u \in \mathcal{D}(x), \alpha u\beta \notin \mathcal{D}(x)$ である．したがって，ある $w \in \mathcal{X}^*$ が存在し， $\alpha uw \in \mathcal{D}(x), \beta \notin \mathcal{P}(w)$ とおける．このとき， $\text{lcp}(u\beta, uw) = |u|$ である．したがって， $u\beta \in \mathcal{P}(r_i)$ だから， $\text{lcp}(r_i, uw) = |u|$ となる．

一方， $uw \in \mathcal{D}(x)$ である．よって，ある整数 k が存在し， $uw \in \mathcal{P}(r_k)$ とおける．このとき， $\text{lcp}(r_i, r_k) = |u|$ となる．

しかしながら， $|v| = |u| + 1$ なので，これは $|v| \neq \text{lcp}(r_i, r_k) + 1$ という仮定に矛盾する．よって， $\mathcal{H}(v) = \phi$ である．

□

補題 1 から $\mathcal{F}(i) \cap \mathcal{G}_H = \phi$ である．すなわち， \mathcal{G}_H の要素 u は $\mathcal{E}(i)$ から探せばよい．ここで， u を効率よく探すため， $\mathcal{E}(i)$ に含まれる 2 つの系列に注目する．すなわち，(2.5) で定義した L 配列を用いて， r_i ごとに，以下の条件を満たす r_i の接頭辞 $\bar{t}_i, \bar{b}_i \in \mathcal{P}(r_i)$ を定義する．

$$|\bar{t}_i| = L[i-1] + 1, \quad (4.7)$$

$$|\bar{b}_i| = L[i] + 1 \leq |r_i|. \quad (4.8)$$

なお， $L[i] = |r_i|$ のとき， \bar{b}_i は存在しないことを注記する．

例として， $x = ababb$ のときの \bar{t}_i, \bar{b}_i を表 4.1 に示す．

表 4.1: $x = ababb$ のときの \bar{t}_i と \bar{b}_i (\bar{b}_i の—は，存在しない事を示す)．

i	r_i	$L[i]$	\bar{t}_i	\bar{b}_i
-1	\perp	-1		
0	λ	0	λ	—
1	ababb	2	a	aba
2	abb	0	abb	a
3	b	1	b	—
4	babb	1	ba	ba
5	bb	-1	bb	λ
6	\perp			

そして， \bar{t}_i と \bar{b}_i を要素に持つ集合 $\mathcal{G}(i)$ を定義する．

$$\mathcal{G}(i) = \begin{cases} \{\bar{t}_i, \bar{b}_i\} & (if \bar{b}_i \text{ exists}) \\ \{\bar{t}_i\} & (otherwise) \end{cases} \quad (4.9)$$

次の補題 2 で， $\mathcal{E}(i)$ と $\mathcal{G}(i)$ の関係を示す．

補題 2 任意の $v \in \mathcal{X}^*$ とある i ($0 \leq i \leq n$) に対して， $v \in \mathcal{E}(i)$ ならば，ある k が存在し，

$$v \in \mathcal{G}(k).$$

証明 最初に，(2.7),(2.8) で定義した被覆集合の最小値と最大値を用いて， $p = \mathcal{C}(v).min, q = \mathcal{C}(v).max$ とおく．

接尾辞配列が整列していることと，(4.6) から， $|v| = \text{lcp}(r_i, r_{p-1}) + 1$ または $|v| = \text{lcp}(r_i, r_{q+1}) + 1$ である．ここで，(2.6) から， $\text{lcp}(r_{p-1}, r_i) = \text{lcp}(r_{p-1}, r_p), \text{lcp}(r_i, r_{q+1}) = \text{lcp}(r_q, r_{q+1})$ である．

L 配列から , $\text{lcp}(r_{p-1}, r_p) + 1 = L[p - 1] + 1, \text{lcp}(r_q, r_{q+1}) + 1 = L[q] + 1$ である .
したがって , $|v| = L[p - 1] + 1$ または $|v| = L[q] + 1$ である . すなわち , $v = \bar{t}_p$ または $v = \bar{b}_q$ となる . よって , (4.9) から , $v \in \mathcal{G}(p)$ または $v \in \mathcal{G}(q)$ である .

以上より , $v \in \mathcal{G}(k)$ である .

□

ここで , 全ての $\mathcal{G}(i)$ の和集合 \mathcal{G} を定義する .

$$\mathcal{G} = \bigcup_{0 \leq i \leq n} \mathcal{G}(i). \quad (4.10)$$

このとき , \mathcal{G}_H は \mathcal{G} に含まれることを , 次の定理 2 で示す .

定理 2 任意の $x \in \mathcal{X}^*$ に対して ,

$$\mathcal{G}_H \subset \mathcal{G}.$$

証明 最初に , $v \in \mathcal{G}_H$ とおくと , (4.4) から , $\mathcal{H}(v) \neq \phi$ である . よって , 補題 1 の対偶から , 任意の i ($0 \leq i \leq n$) について , $v \notin \mathcal{F}(i)$ となる . すなわち ,

$$v \notin \bigcup_{0 \leq i \leq n} \mathcal{F}(i). \quad (4.11)$$

また , $j \in \mathcal{C}(v)$ とおくと , (2.6) から , $v \in \mathcal{P}(r_j)$. さらに , (4.11) から , $v \notin \mathcal{F}(j)$. よって , $v \in \mathcal{E}(j)$ となる . このとき , 補題 2 から , ある k が存在して , $v \in \mathcal{G}(k)$ である . したがって , (4.10) から , $\mathcal{G}_H \subset \mathcal{G}$ である .

□

定理 2 から , \mathcal{G}_H は \mathcal{G} に含まれる . さらに , (4.7), (4.8) から , \mathcal{G} は接尾辞配列と L 配列を走査するだけで容易に求めることができる .

4.4 提案手法の概要

提案手法は , 長さ n の系列 x を入力し , 反辞書 $\mathcal{A}(x)$ を出力する .

反辞書は , \mathcal{G}_H と定理 1 を用いて構築する . ただし , \mathcal{G}_H を求めるのは容易でないことから , 代わりに \mathcal{G} を求める . 定理 2 から , \mathcal{G}_H は \mathcal{G} に含まれている . なお , (4.10) より $|\mathcal{G}|$ は高々 $2(n+1)$ であることから , 全ての $u \in \mathcal{G}$ について \mathcal{G}_H の要素かどうかを調べる . 判別方法は , (4.4) から , $\mathcal{H}(u)$ が空集合かどうかを調べる . この判別方法だと , $\mathcal{H}(u)$ が空集合でないときに , 極小禁止語を容易に求めることができる . すなわち , 定理 1 から , $\mathcal{H}(u)$ の要素を u の先頭に加えると極小禁止語になる .

提案手法 SA2AD の各ステップを以下に示す .

Algorithm SA2AD

```

input   : 入力記号系列  $x$ 
output  : 反辞書  $\mathcal{A}(x)$ 

begin 1
  /* Step 1 : 接尾辞配列の構築 */ 2
   $S \leftarrow \text{construct\_suffix\_array}(x)$  ; 3
  /* Step 2 :  $L$  配列の構築 */ 4
   $L \leftarrow \text{construct\_l\_array}(x, S)$  ; 5
  /* Step 3 :  $\mathcal{G}$  の構築 */ 6
   $\mathcal{G} \leftarrow \text{construct\_G}(S, L)$  ; 7
  /* Step 4 : 反辞書  $\mathcal{A}(x)$  の構築 */ 8
   $\mathcal{A}(x) \leftarrow \text{construct\_antidictionary}(x, S, \mathcal{G})$  ; 9
  return  $\mathcal{A}(x)$  ; 10
end. 11
```

Step 1 と Step 2 では , \mathcal{G} を構築する準備として接尾辞配列 S と L 配列 L を構築する . Step 3 で \mathcal{G} を構築し , Step 4 で $\mathcal{H}(u)$ を求めて $\mathcal{A}(x)$ を構築する .

Step 1 の $\text{construct_suffix_array}(x)$ は , 計算量が線形の既知のアルゴリズム [16] – [20] を用いる . Step 2 の $\text{construct_l_array}(x, S)$ も , 計算量が線形の既知のアルゴリズム [11], [12] を用いる . Step 3 の $\text{construct_G}(S, L)$ の処理内容については 4.6 節で , Step 4 の $\text{construct_antidictionary}(x, S, \mathcal{G})$ については 4.7 節で述べる .

4.5 系列の表現

系列 $u \in \mathcal{G}$ の $\mathcal{H}(u)$ を求めるには , (4.1) より , (2.10) で定義した u の全接頭記号集合 $\mathcal{X}^H(u)$ を求める必要がある .

ここで , (2.7), (2.8) で定義した u の被覆集合の最小値 $\mathcal{C}(u).\text{min}$ と最大値 $\mathcal{C}(u).\text{max}$ を用いると , (2.11) より , $\mathcal{X}^H(u)$ は接尾辞配列 S の $\mathcal{C}(u).\text{min}$ 番目から $\mathcal{C}(u).\text{max}$ 番目を走査して求めることができる . そこで , u は 3 つの組

$$(|u|, \mathcal{C}(u).\text{min}, \mathcal{C}(u).\text{max}) \tag{4.12}$$

で表現する .

ここで , $\mathcal{C}(u).\text{min}, \mathcal{C}(u).\text{max}$ と L 配列との関係を , 次の命題 1 で示す .

命題 1 系列 $u \in \mathcal{D}(x)$ に対して, $p = \mathcal{C}(u).\text{min}, q = \mathcal{C}(u).\text{max}$ とおくと,

$$|u| > L[p-1], |u| > L[q].$$

さらに, $p < q$ ならば,

$$|u| \leq \min\{L[i] \mid p \leq i \leq q-1\}.$$

証明 まず, p のとりうる値 $0 \leq p \leq n$ を, $p = 0$ と $1 \leq p \leq n$ に分けて考える.

はじめに, $p = 0$ のときを考える. (2.5) より $L[p-1] = L[-1] = -1$. 一方, $|u| \geq 0$ である. したがって, $|u| > L[p-1]$ である.

次に, $1 \leq p \leq n$ のときを考える. (2.7) より, $p-1 \notin \mathcal{C}(u)$. すなわち, (2.6) から, $u \notin \mathcal{P}(r_{p-1})$. よって, $|u| > \text{lcp}(r_{p-1}, u)$ である. このとき, 任意の $v \in \mathcal{X}^*$ に対して, $\text{lcp}(r_{p-1}, uv) = \text{lcp}(r_{p-1}, u)$ となる. ところで, $u \in \mathcal{P}(r_p)$ なので, $r_p = uw$ ($w \in \mathcal{X}^*$) と表すと, $L[p-1] = \text{lcp}(r_{p-1}, r_p) = \text{lcp}(r_{p-1}, uw) = \text{lcp}(r_{p-1}, u)$ である. したがって, $|u| > L[p-1]$ である.

また, q のとりうる値 $0 \leq q \leq n$ を, $q = n$ と $0 \leq q \leq n-1$ に分けて考えると, 上と同様の議論から $|u| > L[q]$ が示せる.

さらに, $p < q$ のときを考える. このとき, k ($p \leq k \leq q$) に対して, $u \in \mathcal{P}(r_k)$ である. したがって, i ($p \leq i \leq q-1$) に対して, $|u| \leq \min\{L[i] \mid p \leq i \leq q-1\}$ である.

以上より, 題意が成り立つ.

□

集合 \mathcal{G} の要素のうち, (4.7) で定義した \bar{t}_i ($0 \leq i \leq n$) の長さは定まっているので, 次の系 2 が成り立つ.

系 2 ある i ($0 \leq i \leq n$) が与えられたとき, 系列 $\bar{t}_i \in \mathcal{G}$ に対して,

$$\begin{aligned} \mathcal{C}(\bar{t}_i).\text{min} &= i, \\ \mathcal{C}(\bar{t}_i).\text{max} &= \min\{j \mid |\bar{t}_i| > L[j], i \leq j \leq n\}. \end{aligned}$$

同じく, 集合 \mathcal{G} の要素のうち, (4.8) で定義した $\bar{b}_i \in \mathcal{G}$ の長さも定まっているので, 次の系 3 が成り立つ.

系 3 ある i ($0 \leq i \leq n$) が与えられたとき, 系列 $\bar{b}_i \in \mathcal{G}$ に対して,

$$\begin{aligned} \mathcal{C}(\bar{b}_i).\text{min} &= \max\{j \mid |\bar{b}_i| > L[j-1], 0 \leq j \leq i\}, \\ \mathcal{C}(\bar{b}_i).\text{max} &= i. \end{aligned}$$

提案手法では, \bar{t}_i, \bar{b}_i は (4.12) の3つ組で表現する. なお, (4.7), (4.8) と系 2, 系 3 から, 3つ組の値は i と L だけで求めることができる. 特に, $\mathcal{C}(\bar{t}_i).max$ と $\mathcal{C}(\bar{b}_i).min$ 以外は L を走査することなく決定できる.

4.6 Step 3のアルゴリズム

Step 3 では, 接尾辞配列 S と L 配列 L から, \mathcal{G} を構築する.

集合 \mathcal{G} の要素 \bar{t}_i, \bar{b}_i は (4.12) の3つ組で表現するが, 系 2 と系 3 から, $\mathcal{C}(\cdot).min$ と $\mathcal{C}(\cdot).max$ の求め方が \bar{t}_i と \bar{b}_i とで違う. そこで, \mathcal{G} は2つの集合

$$\mathcal{G}_T = \bigcup_{0 \leq i \leq n} \{\bar{t}_i\}, \quad (4.13)$$

$$\mathcal{G}_B = \bigcup_{0 \leq i \leq n} \{\bar{b}_i\}. \quad (4.14)$$

を求めた後に, 和集合を演算して構築する. なお, \bar{b}_i が存在しない時は, $\{\bar{b}_i\} = \phi$ とする.

さらに, $\mathcal{C}(\cdot).min, \mathcal{C}(\cdot).max$ を求めるときには, (2.11) より, (2.10) で定義した $\mathcal{X}^H(\cdot)$ を用いて S を走査するが, \mathcal{G} の要素が整列していると S の走査を抑えることができる. そこで, \mathcal{G}_T と \mathcal{G}_B が辞書式順序昇順で整列するように構築する. 整列する手法については, この節の最後に述べる. また, \mathcal{G} は, \mathcal{G}_T と \mathcal{G}_B のマージソート [21] で構築する.

そこで, 辞書式順序昇順で整列した後の $\mathcal{G}_T, \mathcal{G}_B, \mathcal{G}$ を

$$\mathcal{G}_T = \{t_0, \dots, t_n\} (t_i \prec t_j, 0 \leq i < j \leq n), \quad (4.15)$$

$$\mathcal{G}_B = \{b_0, \dots, b_k\} (b_i \prec b_j, 0 \leq i < j \leq k) \quad (4.16)$$

$$\mathcal{G} = \{g_0, \dots, g_m\} (g_i \prec g_j, 0 \leq i < j \leq m), \quad (4.17)$$

とあらためて定義する. なお, \bar{b}_i が存在しないことがあるので, $k = |\mathcal{G}_B| - 1 \leq n$ である. \mathcal{G}_T と \mathcal{G}_B が同じ要素を持つことがあるので, $m = |\mathcal{G}| - 1 \leq 2n + 1$ である.

Step 3 のアルゴリズム `construct_G` を以下に示す. アルゴリズム中の `stack` は, S のインデックスを格納するスタックである. `stack.push()` でインデックスを格納し, `stack.pop()` で取り出す. `stack.is_empty()` でスタックが空かどうかを調べ, 空なら `true` を返し, 空でなければ `false` を返す.

Algorithm `construct_G`

input : 接尾辞配列 S , L 配列 L


```

    output : 系列の集合  $\mathcal{G}$ 

begin
    /* Step 3-1 : 初期設定 */
     $j \leftarrow n$ 
    for  $k := 0$  to  $n$  do begin
         $t_k \leftarrow (L[k-1] + 1, k, -1)$ ; /* -1 は未設定を示す */
         $\bar{b}_k \leftarrow (L[k] + 1, -1, k)$ ; /* -1 は未設定を示す */
        if (  $|r_k| \leq L[k]$  ) begin
             $\bar{b}_k \leftarrow \epsilon$ ; /*  $\bar{b}_i$  は存在しない */
             $j \leftarrow j - 1$ ;
        end if;
    end for;
    /* Step 3-2 :  $\mathcal{G}_T$  の構築 */
    stack.clean();
    for  $i := 0$  to  $n$  do begin
        stack.push( $i$ );
        while ( stack.is_empty()  $\neq$  true ) begin
             $k \leftarrow$  stack.pop();
            if (  $|t_k| > L[i]$  )
                 $\mathcal{C}(t_k).max \leftarrow i$ ;
            else begin
                stack.push( $k$ );
                break;
            end if;
        end while;
    end for;
    /* Step 3-3 :  $\mathcal{G}_B$  の構築 */
    stack.clean();
    for  $i := n$  to  $0$  do begin
        if (  $\bar{b}_i \neq \epsilon$  )
            stack.push( $i$ )
        while ( stack.is_empty()  $\neq$  true ) begin
             $k \leftarrow$  stack.pop();
            if (  $|\bar{b}_k| > L[i-1]$  ) begin
                 $\mathcal{C}(\bar{b}_k).min \leftarrow i$ ;

```

```

     $b_j \leftarrow \bar{b}_k;$  35
     $j \leftarrow j - 1;$  36
  else begin 37
    stack.push( $k$ ); 38
    break; 39
  end if; 40
end while; 41
end for; 42
/* Step 3-4 :  $\mathcal{G}$  の構築 */ 43
 $\mathcal{G} \leftarrow \text{merge.sort}(\mathcal{G}_T, \mathcal{G}_B);$  44
return  $\mathcal{G}$ ; 45
end. 46

```

辞書式順序昇順で整列するためには，系列を比較する必要がある．アルゴリズムでは直接扱っていないが，系列の比較は $\mathcal{C}(\cdot).\min$ と $\mathcal{C}(\cdot).\max$ を用いている．この2つの値で系列を比較できることを，次の命題2と命題3で示す．

命題2 任意の $u, v \in \mathcal{D}(x)$ に対して， $\mathcal{C}(u).\min < \mathcal{C}(v).\min$ ならば，

$$u \prec v.$$

証明 最初に， $i = \mathcal{C}(u).\min, j = \mathcal{C}(v).\min$ とおく．(2.6) から， $u \in \mathcal{P}(r_i), v \in \mathcal{P}(r_j)$ である．さらに，(2.6) と $i < j$ から， $i \notin \mathcal{C}(v)$ ．したがって， $u \notin \mathcal{P}(r_j)$ となる．

すなわち， $r_i \prec r_j$ なので， $u \prec v$ である．

□

命題3 任意の $u, v \in \mathcal{D}(x)$ に対して， $\mathcal{C}(u).\min = \mathcal{C}(v).\min, \mathcal{C}(u).\max > \mathcal{C}(v).\max$ ならば，

$$u \prec v.$$

証明 最初に，(2.9) から， $\mathcal{C}(v) \subsetneq \mathcal{C}(u)$ である．ここで，背理法を使って証明する．いま， $|u| \geq |v|$ と仮定する．このとき，ある整数 i が存在し， $i \in \mathcal{C}(u) \cap \mathcal{C}(v)$ とおける．すなわち， $u \in \mathcal{P}(r_i), v \in \mathcal{P}(r_i)$ である．よって， $|u| \geq |v|$ から， $\mathcal{C}(v) \subset \mathcal{C}(u)$ となる．しかしながら，これは， $\mathcal{C}(v) \subsetneq \mathcal{C}(u)$ に矛盾する．したがって， $|u| < |v|$ である．

そして, $|u| < |v|$ と $u \in \mathcal{P}(r_i), v \in \mathcal{P}(r_i)$ から, $u \in \mathcal{P}(v)$ となる.
 以上より, $u \prec v$ である.

□

集合 \mathcal{G}_T の構築は, Step 3-1 と 3-2 で行われる. このステップにて, \bar{t}_k を用いずに t_k を直接求めているのは, 系 2 と命題 2 から, $\bar{t}_0 \prec \bar{t}_1 \prec \dots \prec \bar{t}_n$ だからである. すなわち, $t_k = \bar{t}_k$ である. (4.12) の 3 要素のうち, Step 3-1 では, (4.7) から $|t_k|$ を, 系 2 から $\mathcal{C}(t_k).min$ を求めている.

そして, Step 3-2 で, $\mathcal{C}(t_k).max$ を求めている. $\mathcal{C}(t_k).max$ は系 2 で与えられる. すなわち, $\mathcal{C}(t_k).max$ を求めるには, k 以上で最初に L の値が $|t_k|$ より小さくなるインデックスを探せばよい. これは, 直近上位要素発見問題 (Nearest Larger Neighbor Problem)[22] である.

そこで, Finding Nearest Larger Neighborhood Algorithm[22] を Step 3-2 の 13 行目から 25 行目で用いている. このアルゴリズムではスタックを活用し, i を 0 から n まで増やしながらか $|t_k| = i$ となる k の値を探している. これにより, 1 回の走査で全ての t_k を求めることができる.

集合 \mathcal{G}_B の構築は, Step 3-1 と 3-3 で行われる. Step 3-1 では, (4.8) から $|b_k|$ を, 系 3 から $\mathcal{C}(b_k).min$ を求めている. そして, Step 3-3 で, $\mathcal{C}(\bar{b}_k).min$ を求めている. $\mathcal{C}(\bar{b}_k).min$ は系 3 で与えられる. これも, k 以下で最初に L の値が $|b_k|$ より小さくなるインデックスを探す直近上位要素発見問題である. そこで, i を n から 0 まで減らしながらか L を走査する Finding Nearest Larger Neighborhood Algorithm を Step 3-3 の 27 行目から 42 行目で用いている.

ところで, i を n から 0 に向かって走査するため, $\mathcal{C}(\bar{b}_k).min$ の値が大きい順に求めることになる. しかも, $\mathcal{C}(\bar{b}_k).min$ が同じ値の場合は, $\mathcal{C}(\bar{b}_k).max$ の値が小さい順に求める. すなわち, 命題 2 と命題 3 から, $\mathcal{C}(\bar{b}_k).min$ の値を求める順番は, \bar{b}_k の辞書式順序降順である. この性質を利用するために, Finding Nearest Larger Neighborhood Algorithm に 35, 36 行目の処理を加え, 辞書式順序昇順に整列した b_k を求めている.

4.7 Step 4 のアルゴリズム

Step 4 では, Step 3 で求めた集合 \mathcal{G} の全ての要素 g_i ($0 \leq i \leq |\mathcal{G}| - 1$) に対して, (4.1) で定義した $\mathcal{H}(g_i)$ を求めて, 反辞書を構築する. (4.1) より, $g_i \neq \lambda$ の場合に $\mathcal{H}(g_i)$ を求めるには, (2.10) で定義した全接頭記号集合 $\mathcal{X}^H(\pi(g_i))$ と

$\mathcal{X}^H(g_i)$ が必要である．ここで，

$$l(i) = \max\{j | j < i, |g_j| < |g_i|\} \quad (4.18)$$

を定義し，次の補題 3 と命題 4 を示す．

補題 3 任意の g_i ($\neq \lambda$) に対して，

$$g_{l(i)} \in \mathcal{P}(g_i).$$

証明 背理法を使って証明する．いま， $g_{l(i)} \notin \mathcal{P}(g_i)$ と仮定する．このとき，ある記号 $\alpha, \beta \in \mathcal{X}$ ($\alpha < \beta$) と系列 $u, v, w \in \mathcal{X}^*$ ($|u| < |v|$) が存在し， $g_{l(i)} = w\alpha u, g_i = w\beta v$ とおける．よって，(4.6) と補題 2 から，ある整数 j が存在し， $g_j = w\beta$. このとき， $|g_j| < |g_i|, g_{l(i)} \prec g_j \prec g_i$ から， $l(i) < j < i$ である．しかしながら，これは，(4.18) に矛盾する．よって， $g_{l(i)} \in \mathcal{P}(g_i)$ である．

□

命題 4 任意の g_i ($\neq \lambda$) に対して，

$$\mathcal{X}^H(\pi(g_i)) = \mathcal{X}^H(g_{l(i)}).$$

証明 表記を簡潔にするために， $u = g_{l(i)}, v = \pi(g_i)$ とおく．(4.18) から， $|u| \leq |v|$ である．そこで， $|u| = |v|$ と $|u| < |v|$ の 2 つの場合に分けて証明する．

最初に， $|u| = |v|$ の場合を証明する．補題 3 と $v \in \mathcal{P}(g_i)$ から， $u = v$. したがって， $\mathcal{X}^H(u) = \mathcal{X}^H(v)$ である．

次に， $|u| < |v|$ の場合を証明する．

補題 3 と (4.18) から， u は $\mathcal{P}(g_i) \cap \{g_0, g_1, \dots, g_{i-1}\}$ の最長系列である．

ここで，背理法を使って証明する．いま， $\mathcal{C}(u) \setminus \mathcal{C}(v) \neq \phi$ と仮定する．このとき，ある整数 j, k が存在し， $j \in \mathcal{C}(u) \cap \mathcal{C}(v), k \in \mathcal{C}(u) \setminus \mathcal{C}(v)$ とおける．すなわち，(2.6) と補題 3 から， $|u| \leq \text{lcp}(r_j, r_k) \leq |v| - 1$ となる．したがって，(4.6) から，ある $w \in \mathcal{X}^*$ が存在し， $w \in \mathcal{E}(j), |w| = \text{lcp}(r_j, r_k) + 1$ とおける．このとき，補題 2 から， $w \in \mathcal{G}$ である．ところで， $|u| < |w| < |g_i|$ と $u, w \in \mathcal{P}(g_i)$ から， $u \prec w \prec g_i$ である．よって， $w \in \{g_0, g_1, \dots, g_{i-1}\}$. しかしながら，これは， u が $\mathcal{P}(g_i) \cap \{g_0, g_1, \dots, g_{i-1}\}$ の最長系列であることに矛盾する．したがって， $\mathcal{C}(u) \setminus \mathcal{C}(v) = \phi$ である．すなわち，(2.11) から， $\mathcal{X}^H(u) = \mathcal{X}^H(v)$ である．

以上より，題意が成り立つ．

□

系 4 任意の g_i ($\neq \lambda$) に対して ,

$$\mathcal{H}(g_i) = \mathcal{X}^H(g_{l(i)}) \setminus \mathcal{X}^H(g_i).$$

Step 4 のアルゴリズム `construct_antidictionary` を以下に示す . アルゴリズム中の $\Lambda[\cdot]$ は , x のインデックスを格納する長さ $|\mathcal{X}|$ の配列である . `stack` は , g_i の添え字を格納するスタックである .

なお , $\mathcal{X}^H(i)$ に記号 α を追加する処理を $\mathcal{X}^H(i).\text{add}(\alpha)$ と記す . 同様にして , $\mathcal{H}(u)$ に記号 α を追加する処理を $\mathcal{H}(u).\text{add}(\alpha)$ と記す . $\mathcal{A}(x)$ に系列 u を追加する処理を $\mathcal{A}(x).\text{add}(u)$ と記す .

Algorithm `construct_antidictionary`

input : 入力系列 x , 接尾辞配列 S , 系列の集合 \mathcal{G}

output : 反辞書 $\mathcal{A}(x)$

```

begin                                                                    1
  /* Step 4-1 : 初期設定 */                                              2
  for (  $\forall \alpha \in \mathcal{X}$  ) do                                       3
     $\Lambda[\alpha] \leftarrow n + 1$ ;                                       4
  /* Step 4-2 :  $\mathcal{X}^H(g_i)$  の構築 */                                  5
   $i \leftarrow |\mathcal{G}| - 1$                                            6
  for  $k := n$  to 0 do begin                                              7
    if (  $S[k] - 1 \geq 1$  )                                              8
       $\Lambda[x_{S[k]-1}] \leftarrow k$ ; /* 最終インデックスの更新 */      9
    while (  $i \geq 0$  and  $\mathcal{C}(g_i).\text{min} = k$  ) begin                  10
      for (  $\forall \alpha \in \mathcal{X}$  ) do begin                               11
        if (  $\Lambda[\alpha] \leq \mathcal{C}(g_i).\text{max}$  )                     12
           $\mathcal{X}^H(g_i).\text{add}(\alpha)$ ;                               13
        end for;                                                         14
       $i \leftarrow i - 1$ ;                                              15
    end while;                                                            16
  end for;                                                                17
  /* Step 4-3 :  $\mathcal{H}(g_i)$  の構築 */                                      18
  /* Step 4-3.1 :  $\mathcal{H}(g_0)$  の設定 */                                  19
  for (  $\forall \alpha \in \mathcal{X}$  ) do begin                                   20
    if (  $\Lambda[\alpha] = n + 1$  )                                       21
       $\mathcal{H}(g_0).\text{add}(\alpha)$ ; /*  $\mathcal{H}(\lambda)$  */                22

```

```

end for; 23
/* Step 4-3.2 :  $l(i)$  の設定 */ 24
 $p \leftarrow 0$  25
stack.clean(); 26
for  $i := 1$  to  $|\mathcal{G}| - 1$  do begin 27
    while (stack.is_empty()  $\neq$  true and  $|g_p| \geq |g_i|$ ) 28
         $p \leftarrow$  stack.pop(); 29
     $l(i) \leftarrow p$ ; 30
    stack.push( $p$ ); 31
     $p \leftarrow i$ ; 32
end for; 33
/* Step 4-3.3 :  $\mathcal{H}(g_i)$  ( $i \geq 1$ ) の設定 */ 34
for  $i := 1$  to  $|\mathcal{G}| - 1$  do 35
     $\mathcal{H}(g_i) \leftarrow \mathcal{X}^H(g_{l(i)}) \setminus \mathcal{X}^H(g_i)$ ; 36
/* Step 4-4 :  $\mathcal{A}(x)$  の構築 */ 37
for  $i := 0$  to  $|\mathcal{G}| - 1$  do 38
    for ( $\forall \alpha \in \mathcal{H}(g_i)$ ) do 39
         $\mathcal{A}(x)$ .add( $\alpha g_i$ ); 40
return  $\mathcal{A}(x)$ ; 41
end. 42

```

Step 4-2 では, k を n から 0 まで減らしながら, $\mathcal{X}^H(g_i)$ ($0 \leq i \leq |\mathcal{G}| - 1$) を構築している. k が減るたびに記号 $x_{\mathbb{S}[k]-1}$ を調べ, $\Lambda[x_{\mathbb{S}[k]-1}]$ に k を格納する. すなわち, k の値を固定したとき,

$$\Lambda[\alpha] = \min\{j | x_{\mathbb{S}[j]-1} = \alpha, k \leq j \leq n\} \quad (\alpha \in \mathcal{X})$$

である. そして, (2.11) から, $\mathcal{C}(g_i).$ min = k となる全ての g_i について, それぞれ $k \leq \Lambda[\alpha] \leq \mathcal{C}(g_i).$ max を満たす α を調べて, $\mathcal{X}^H(g_i)$ を求める.

Step 4-3.1 では, $\mathcal{H}(g_0)$, すなわち $\mathcal{H}(\lambda)$ を求めている. これは, (4.1) から, 未出現記号の集合である. そこで, $\Lambda[\alpha]$ を調べて求める. また, Step 4-3.2 では (4.18) の $l(i)$ を求め, Step 4-3.3 で $\mathcal{H}(g_i)$ を求める.

最後に Step 4-4 で, 定理 1 から, $\mathcal{H}(g_i)$ と g_i を用いて全ての極小禁止語を求める, 反辞書 $\mathcal{A}(x)$ を構築する.

4.8 計算時間の評価

提案手法 SA2AD は , Step 1 から Step 4 までの 4 ステップで構成されている . このうち Step 1 と Step 2 は , 計算時間と記憶量が線形の既知の手法 [16] – [20], [11], [12] を用いる . Step 1 と Step 2 の計算時間をそれぞれ $T_1(n), T_2(n)$ とおくと , $T_1(n) = O(n), T_2(n) = O(n)$ である .

ここで , 非負整数上に定義された実関数 $f(n), g(n)$ に対して , $f(n) = O(g(n))$ であるとは , ある定数 $E > 0$ とある $n_0 \in \mathbb{N}$ が任意の $n > n_0$ に対して ,

$$|f(n)| \leq E|g(n)|$$

が成立することとする .

Step 3 の計算時間を $T_3(n)$ とおくと , 次の補題 4 を示せる .

補題 4 入力系列 x の長さ n に対して ,

$$T_3(n) = O(n) .$$

証明 Step 3 のアルゴリズム `construct_G` の Step 3-1, 3-2, 3-3, 3-4 の計算時間をそれぞれ $T_{31}, T_{32}, T_{33}, T_{34}$ とおく . このうち , $T_{31} = O(n)$ は自明である .

Step 3-2 の 16-24 行目にある while ループは , `stack` が空でないときに繰り返す . また , ループ中に `stack.pop()` を必ず行う . したがって , Step 3-2 の処理が終わるまでに while ループを繰り返す総回数の上限は , Step 3-2 の処理が終わるまでに `stack.push(·)` した総回数に等しい . Step 3-2 は , 全体が for ループとなっており , 1 回のループにつき最大で 2 回 `stack.push(·)` する . したがって , while ループを繰り返す総回数の上限は $2n + 2$. よって , $T_{32} = O(n)$.

Step 3-3 も同様に , 31-41 行目の while ループを繰り返す総回数の上限は $2n + 2$. よって , $T_{33} = O(n)$.

Step 3-4 では , \mathcal{G}_T と \mathcal{G}_B のマージソートを行っている . このステップまでに $\mathcal{G}_T, \mathcal{G}_B$ はソート済みである . (4.13), (4.14) から , $|\mathcal{G}_T| \leq n + 1, |\mathcal{G}_B| \leq n + 1$. よって , $T_{34} = O(n)$.

以上より , $T_3(n) = T_{31} + T_{32} + T_{33} + T_{34} = O(n)$ である .

□

Step 4 の計算時間を $T_4(n)$ とおくと , 次の補題 5 を示せる .

補題 5 入力系列 x の長さ n に対して ,

$$T_4(n) = O(n) .$$

証明 Step 4 のアルゴリズム `construct_antidictionary` の Step 4-1, 4-2, 4-3, 4-4 の計算時間をそれぞれ $T_{41}, T_{42}, T_{43}, T_{44}$ とおく . このうち , $T_{41} = O(1)$ は自明である . (4.9), (4.10) から , $|\mathcal{G}| \leq 2n + 2$. よって , $T_{44} = O(n)$.

Step 4-2 の 10-16 行目にある while ループは , i が負のときは繰り返さない . また , i の初期値は $|\mathcal{G}| - 1$ で , ループ中に必ず i の値を -1 する . $|\mathcal{G}| \leq 2n + 2$ から , Step 4-2 の処理が終わるまでに while ループを繰り返す総回数の上限は , $2n + 2$. よって , $T_{42} = O(n)$.

次に , Step 4-3.1, 4-3.2, 4-3.3 の計算時間をそれぞれ $T_{43.1}, T_{43.2}, T_{43.3}$ とおく . このうち , $T_{43.1} = O(1)$ は自明である . さらに , $|\mathcal{G}| \leq 2n + 2$ から , $T_{43.3} = O(n)$. Step 4-3.2 の 28-29 行目にある while ループは , `stack` が空のときは繰り返さない . また , ループ中に必ず `stack.pop()` する . したがって , Step 4-3.2 の処理が終わるまでに while ループを繰り返す回数の上限は , Step 4-3.2 の処理が終わるまでに `stack.push()` した総回数に等しい . Step 3-2 は , 全体が for ループとなっており , 1 回のループにつき 1 回 `stack.push()` する . したがって , while ループを繰り返す総回数の上限は $n + 1$. よって , $T_{43.2} = O(n)$. すなわち , $T_{43} = T_{43.1} + T_{43.2} + T_{43.3}$ から , $T_{43} = O(n)$.

以上より , $T_4(n) = T_{41} + T_{42} + T_{43} + T_{44} = O(n)$ である .

□

提案アルゴリズム SA2AD の計算時間を $T(n)$ とおくと , 次の定理 3 を示せる .

定理 3 入力系列 x の長さ n に対して ,

$$T(n) = O(n) .$$

証明 補題 4 , 補題 5 と , $T(n) = T_1(n) + T_2(n) + T_3(n) + T_4(n)$ から , $T(n) = O(n)$ である .

□

4.9 記憶量の評価

アルゴリズム SA2AD の Step 1, Step 2, Step 3, Step 4 のそれぞれの記憶量を $M_1(n), M_2(n), M_3(n), M_4(n)$ とおく . なお , 記号とインデックス , そして系列長は固定幅の記憶量で格納するので , 記憶量を定数 c とする . なお , 系列は (4.12) で表現しているので , 記憶量を定数 $3c$ とする .

Step 1 と Step 2 は , 計算時間と記憶量が $O(n)$ の既知の手法 [16] – [20], [11], [12] を用いるから , $M_1(n) = O(n), M_2(n) = O(n)$ である .

Step 3 のアルゴリズム `construct_G` の記憶量を , 次の補題 6 で示す .

補題 6 入力系列 x の長さ n に対して ,

$$M_3(n) = O(n) .$$

証明 全ての $t_i \in \mathcal{G}_T$ を格納する記憶量は , (4.15) から , $3(n+1)c$. 全ての $b_i, \bar{b}_i \in \mathcal{G}_B$ を格納する記憶量は , (4.14), (4.16) から , 高々 $6(n+1)c$. そして , `stack` には , 接尾辞 S のインデックスを重複することなく格納するから , 記憶量は高々 $(n+1)c$. また , (4.10) から , $|\mathcal{G}| \leq 2(n+1)$. よって , \mathcal{G} を格納する記憶量は高々 $6(n+1)c$.
以上より , $M_3(n) \leq 16(n+1)c$. よって , $M_3(n) = O(n)$ である .

□

Step 4 のアルゴリズム `construct_antidictionary` の記憶量を , 次の補題 7 で示す .

補題 7 入力系列 x の長さ n に対して ,

$$M_4(n) = O(n) .$$

証明 配列 Λ は , 長さ $|\mathcal{X}|$ でインデックスを格納するから , 記憶量は $|\mathcal{X}|c$. そして , \mathcal{G} を格納する記憶量は , (4.10) から , 高々 $6(n+1)c$.

また , $\mathcal{X}^H(\cdot)$ は高々 $|\mathcal{X}|$ 個の記号からなる集合である . すなわち , $g_i \in \mathcal{G}$ ごとに $\mathcal{X}^H(g_i)$ を格納するので , 記憶量は高々 $6|\mathcal{X}|(n+1)c$. なお , $\mathcal{H}(\cdot)$ も高々 $|\mathcal{X}|$ 個の記号からなる集合である . そして , g_i ごとに格納するから , 記憶量は高々 $6|\mathcal{X}|(n+1)c$.

つぎに , $l(i)$ は , 長さ $|\mathcal{G}|$ の配列でインデックスを格納するから , 記憶量は高々 $2(n+1)c$. さらに , `stack` は , インデックスを格納し , 最大で $|\mathcal{G}|$ 回の `stack.push()` を実行するから , 記憶量は高々 $2(n+1)c$. また , $\mathcal{A}(x)$ は , 記号と系列を格納し , 最大で $|\mathcal{X}||\mathcal{G}|$ 個の要素を加えるから , 記憶量は高々 $8|\mathcal{X}|(n+1)c$.

以上より , $M_4(n) \leq (20|\mathcal{X}| + 10)(n+1)c + |\mathcal{X}|c$. よって , $M_4(n) = O(n)$ である .

□

提案アルゴリズム SA2AD の記憶量を $M(n)$ とおくと , 次の定理 4 を示せる .

定理 4 入力系列 x の長さ n に対して ,

$$M(n) = O(n) .$$

証明 補題 6 , 補題 7 と , $M(n) = M_1(n) + M_2(n) + M_3(n) + M_4(n)$ から , $M(n) = O(n)$ である .

□

4.10 計算機実験

提案手法と接尾辞木を用いた手法とを、計算機実験で比較した。計算機のCPUはCore2 Duo E8500 (3.16GHz, 2コア)、メモリ量は4GByte、OSはWindows XP SP3 32bit版である。プログラムは両方ともC++言語で作成した。

入力系列は、 $|\mathcal{X}| = 256$ の独立同一分布をもつ情報源で、各記号の生起確率は等確率である。入力系列長は、 50×10^3 ごとに用意した。それぞれ計算時間と記憶量を3回測定し、平均を求めた。なお、OSの制限でアプリケーションが使用できる最大メモリ量は2Gbyteとなっている。そのため、接尾辞木を用いた手法では入力系列長 200×10^3 で、提案手法では入力系列長 450×10^3 で、使用メモリ量が多すぎて実行不可能となった。

計算時間の比較結果を図4.1に、記憶量の比較結果を図4.2に示す。

提案手法の計算時間と記憶量の両方とも、 $O(n)$ であることが計算機実験でも示された。また、接尾辞木を用いた手法と比べて、提案手法の計算量は、計算時間が約20分の1に、記憶量が約2.5分の1に改善された。

さらに、記憶量の内訳について調べた。

接尾辞木を用いた手法では、記憶量の約99%がMF-linkの管理に使われていた。MF-linkとは、接尾辞木の中間ノードのうち、枝を2本以上持つ中間ノードから伸びる特別なリンクである。接尾辞木にMF-linkを追加することで反辞書を構築する。MF-linkを追加した中間ノードからは、高々 $|\mathcal{X}|$ 本のMF-linkが伸びているため、MF-linkの管理に使用する記憶量は $|\mathcal{X}|n$ に比例する。今回の計算機実験において、MF-linkの管理に使用している記憶量は約 $40|\mathcal{X}|n = 10000n$ byteであった。

一方、提案手法では、記憶量の約94~95%が $\mathcal{A}(x)$ の記憶に使われていた。 $\mathcal{A}(x)$ に格納する極小禁止語は、先頭記号と後続する系列に分けて表現している。すなわち、 $\mathcal{A}(x)$ に使用する記憶量も $|\mathcal{X}|n$ に比例する。今回の計算機実験において、 $\mathcal{A}(x)$ に使用している記憶量は、約 $14.8|\mathcal{X}|n \sim 16.4|\mathcal{X}|n = 3800n \sim 4200n$ byteとなった。すなわち、提案方式で改善された記憶量は、MF-linkの管理に使用している記憶量と $\mathcal{A}(x)$ の記憶量との差であることがわかった。

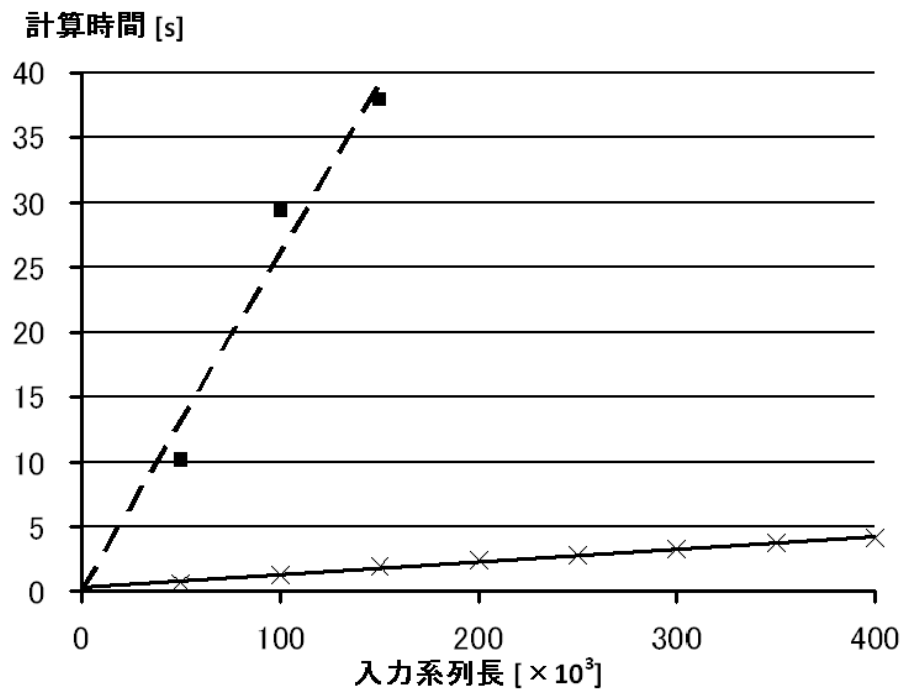


図 4.1: 計算時間の比較 (実線: 提案手法, 点線: 接尾辞木を用いた手法)

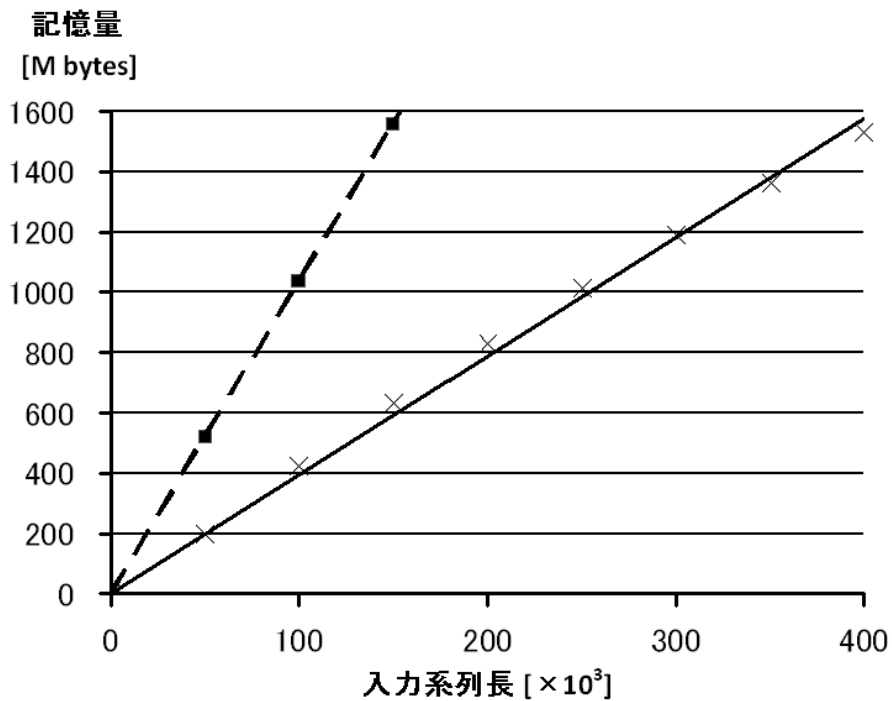


図 4.2: 記憶量の比較 (実線: 提案手法, 点線: 接尾辞木を用いた手法)

次に，アルファベットサイズ m ($= |\mathcal{X}|$) を変化させたときの計算時間と記憶量の変化を測定した．用意したのは， m がそれぞれ 2, 50, 100, 150, 200, 250 の独立同一分布をもつ情報源で，各記号の生起確率は等確率である．入力系列長は， 50×10^3 ごとに，系列長 $50 \times 10^3 \sim 250 \times 10^3$ まで用意した．それぞれ計算時間と記憶量を 100 回測定し，平均を求めた．

計算時間の測定結果を図 4.3 に，記憶量の測定結果を図 4.4 に示す．

いずれのアルファベットサイズでも，計算時間と記憶量は入力系列長に対して線形であることが実験結果でも示せた．また，計算時間と記憶量はアルファベットサイズにも比例することがわかる．これは，Step 4 のアルゴリズム `construct_antidictionary` の 3～4 行目，11～14 行目，20～23 行目に m 回繰り返す処理があるためである．さらに，36 行目の $\mathcal{H}(g_i)$ を求める処理と 39 行目の `for` 文の実装においても，最大で m 回繰り返す処理を用いている．

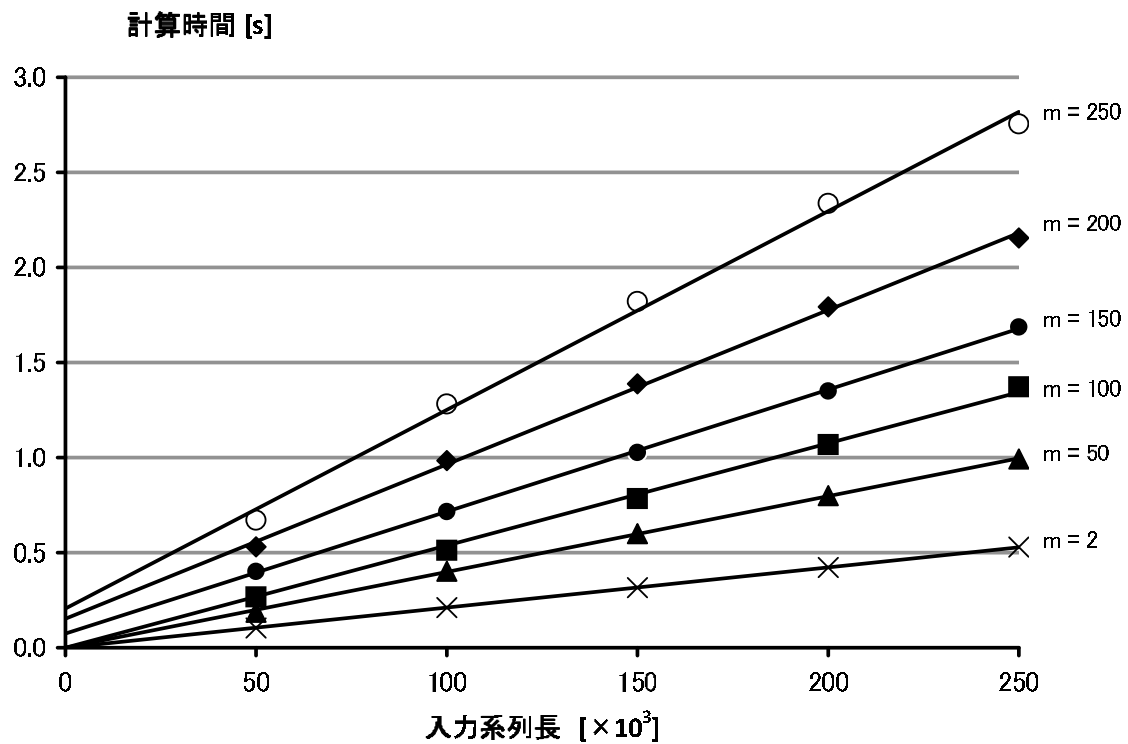


図 4.3: アルファベットサイズ m ごとの計算時間

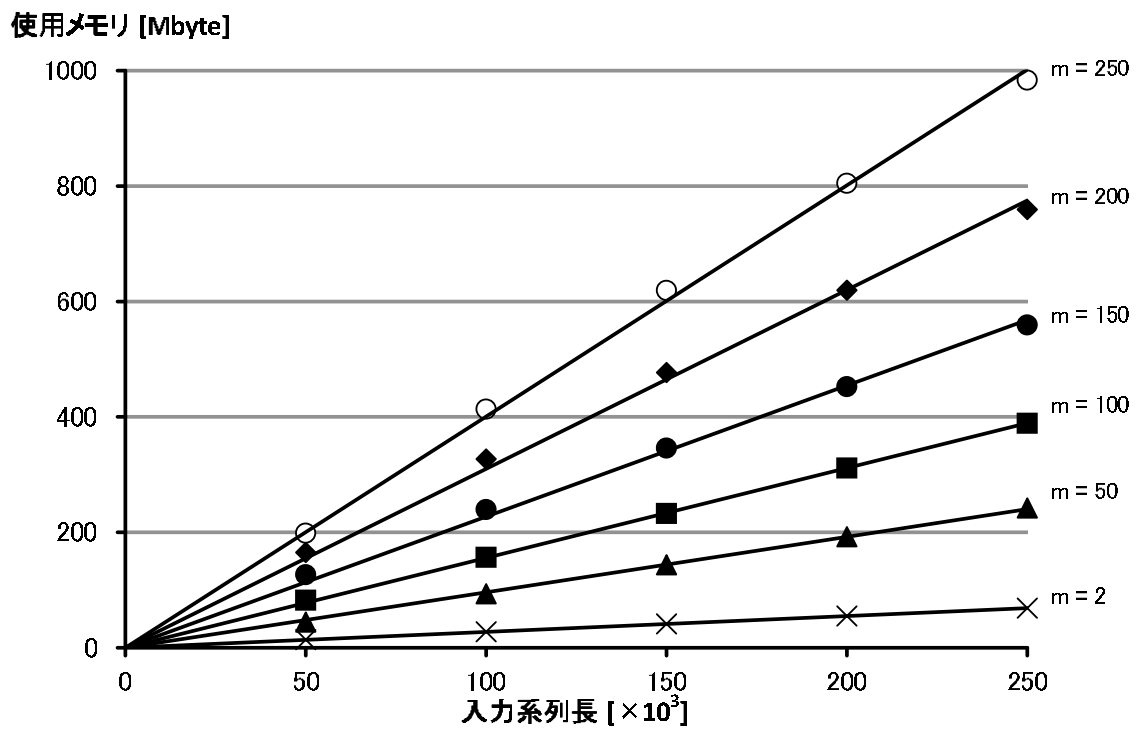


図 4.4: アルファベットサイズ m ごとの記憶量

第5章 動的反辞書構築の提案

5.1 静的反辞書構築と動的反辞書構築

静的反辞書構築とは，入力系列を一度すべて読み込み終わってから反辞書を構築することである．第4章で提案した手法も静的反辞書構築である．接尾辞配列やL配列のように，既存のデータ構造の多くは入力系列をすべて読み込み終わった状態から構築するアルゴリズムが多い．そのため，反辞書を静的に構築する場合には，これらのデータ構造を活用できる利点がある．

ただし，反辞書を使って入力系列自身を圧縮したり，調べたりするときには，静的反辞書構築ではオーバーヘッドを生じることになる．例えば，データ圧縮の反辞書法 [5] では，圧縮系列に加えて，データ圧縮で使う反辞書オートマトンをオーバーヘッドとして復号側へ送る．不整脈検出 [8] では，学習系列を準備して反辞書オートマトンを作成する．

一方，動的反辞書構築とは，入力系列を1記号ずつ読み込みながら逐次的に反辞書を構築することである．入力記号に合わせながら反辞書を変化させていくので，上記のオーバーヘッドを無くすることができる．また，1記号読み込んだ後の反辞書の変化が少ない場合，第4章の提案手法よりも効率がよくなる可能性がある．

この章では，記号の入力ごとに，逐次的に反辞書オートマトンを変化させていく新しい手法を提案する．反辞書オートマトンについては，5.3節に後述する．

5.2 表記の拡張

入力系列

$$\mathbf{x}_1^n = x_1 x_2 \cdots x_n$$

は，記号 x_{n+1} の入力によって系列の末尾に記号 x_{n+1} が加わった系列

$$\mathbf{x}_1^{n+1} = x_1 x_2 \cdots x_n x_{n+1}$$

に変化する．本章までは静的反辞書構築を扱っていたので，便宜上， n を意識しない形で表記していた．しかし，この章からは記号 x_{n+1} を入力する動的反辞書構築なので，以下のように n を明示する形に表記を拡張する．

- 系列 x_1^n を， x^n と表記する．
- 第2章の (2.10) で定義した，系列 x^n 上の u の全接頭記号集合 $\mathcal{X}^H(u)$ を， $\mathcal{X}_n^H(u)$ と表記する．
- 第2章の (2.12) で定義した，系列 x^n 上の u の全接尾記号集合 $\mathcal{X}^T(u)$ を， $\mathcal{X}_n^T(u)$ と表記する．
- 第4章の (4.1) で定義した，系列 x^n 上の $\mathcal{H}(u)$ を， $\mathcal{H}_n(u)$ と表記する．
- 第4章の (4.3) で定義した，系列 x^n 上の $\mathcal{T}(u)$ を， $\mathcal{T}_n(u)$ と表記する．

5.3 反辞書オートマトン

反辞書オートマトンとは，極小禁止語を受理する決定性有限オートマトンである．反辞書法 [5] で用いられている．

図 5.1 に $x^{10} = \text{abcaababac}$ から構築した反辞書オートマトンを示す．このとき，

$$\mathcal{A}(x^{10}) = \{\text{aaa}, \text{aabac}, \text{aabc}, \text{aac}, \text{aca}, \text{baa}, \text{babab}, \text{babac}, \text{bb}, \text{cab}, \text{cac}, \text{cb}, \text{cc}\}$$

である．

一番左上のノードが開始状態である．本論文では，ルートノードと呼ぶことにする．黒色のノードが極小禁止語を受理する状態である．本論文では，MFW ノードと呼ぶことにする．それ以外のノードを中間ノードと呼ぶことにする．

また，任意のノード p に対して，ルートノードから最短の遷移で到達する経路を最短パスと呼ぶことにする．さらに，最短パス上にあるラベルをルートノードから順番に接続した系列を p のラベル列と呼び， $L(p)$ と表記する．MFW ノードのラベル列は，受理する極小禁止語を表す．

この反辞書オートマトンには，任意の系列 $v \in \mathcal{X}^*$ を入力することができる．ルートノードから開始し，入力系列の先頭から順番に，入力記号と同じラベルを持つ遷移を進んでいく．極小禁止語を受理するか，すべての入力記号の

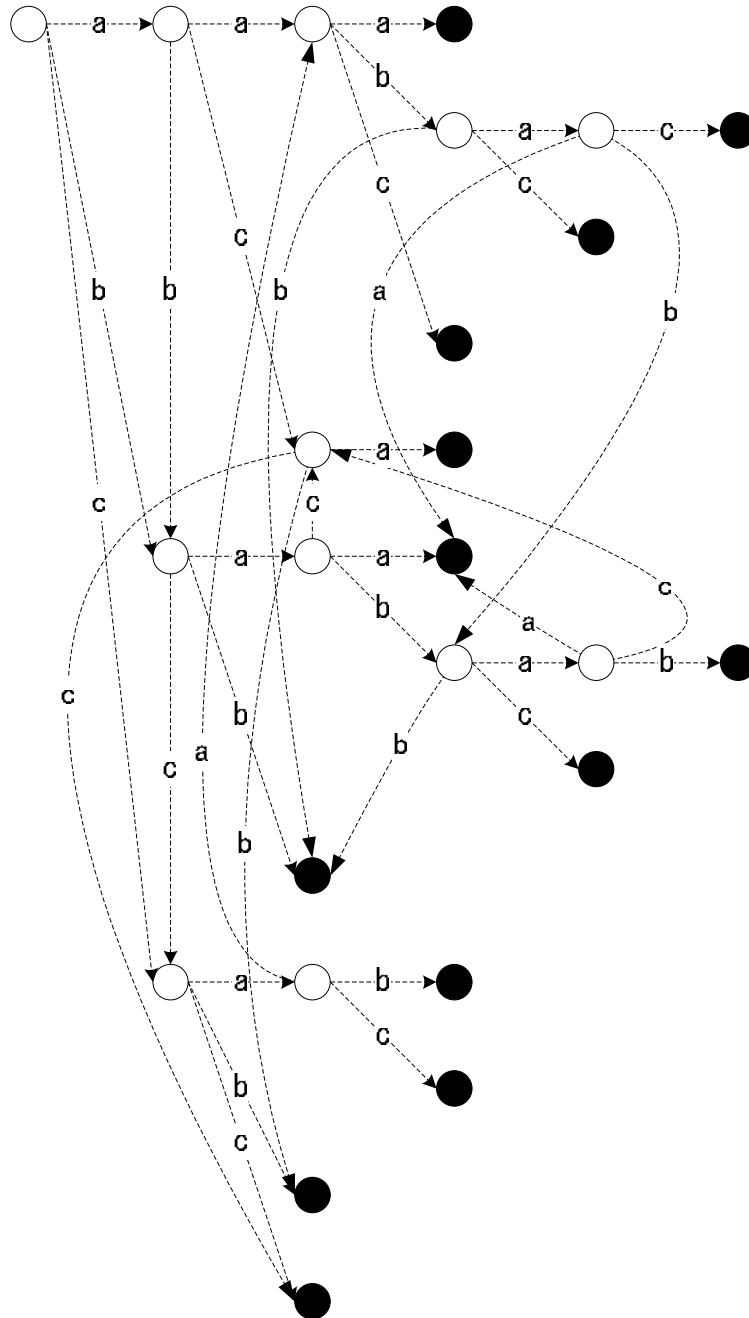


図 5.1: $x^{10} = abcaababac$ の反辞書オートマトン．黒丸のノードが極小禁止語を示す．

遷移を終了すると遷移は終了である．このとき到達したノードを v のローカスと呼び， $\mathcal{N}(v)$ と表記する．

反辞書オートマトンは，反辞書トライを拡張して作ることができる．すなわち，反辞書トライの葉ノード以外のノード q に，すべての記号 $\beta \in \mathcal{X}$ のラベルを持つ遷移を追加する．このとき， $L(q)\beta$ をラベル列に持つノードが存在しない時は $\sigma(L(q)\beta)$ をラベル列に持つノードを探して遷移する．このノードも存在しないときには，さらに $\sigma^2(L(q)\beta)$ をラベル列に持つノードを探して遷移する．以降，遷移先のノードが見つかるまでラベル列の先頭記号を削っていきながら遷移先のノードを探していく．

言いかえると， $L(q)\beta$ の接尾辞をラベル列にもつノードの中で，ラベル列の長さが最大のノードへの遷移を，反辞書トライに追加する．

なお，図 5.1 には書かれていないが，本論文では全てのノード p には， $L(p)$ の接尾辞をラベル列にもつノードの中で，ラベル列の長さが最大のノードへのリンクを持っている．このリンクを suffix link と呼ぶことにする．ノード p の suffix link 先のノードを $s(p)$ と表記する．

5.4 既知の結果

入力系列 x^n ($n \geq 1$) の末尾に記号 x_{n+1} を加えて x^{n+1} になったとき，反辞書は $\mathcal{A}(x^n)$ から $\mathcal{A}(x^{n+1})$ へ変化する．このとき， $\mathcal{A}(x^n)$ から削除された極小禁止語を $\mathcal{A}_{\text{DEL}}(x^n \rightarrow x^{n+1})$ とおき，新しく加わった極小禁止語を $\mathcal{A}_{\text{NEW}}(x^n \rightarrow x^{n+1})$ とおく．すなわち，

$$\mathcal{A}_{\text{DEL}}(x^n \rightarrow x^{n+1}) = \{w \mid w \in \mathcal{A}(x^n), w \notin \mathcal{A}(x^{n+1})\}, \quad (5.1)$$

$$\mathcal{A}_{\text{NEW}}(x^n \rightarrow x^{n+1}) = \{w \mid w \notin \mathcal{A}(x^n), w \in \mathcal{A}(x^{n+1})\} \quad (5.2)$$

である．明らかに， $\mathcal{A}_{\text{DEL}}(x^n \rightarrow x^{n+1})$ と $\mathcal{A}_{\text{NEW}}(x^n \rightarrow x^{n+1})$ は互いに素である．また， $\mathcal{A}(x^n)$ と $\mathcal{A}_{\text{NEW}}(x^n \rightarrow x^{n+1})$ も互いに素である．そして，

$$\mathcal{A}(x^{n+1}) = \{\mathcal{A}(x^n) \setminus \mathcal{A}_{\text{DEL}}(x^n \rightarrow x^{n+1})\} \cup \mathcal{A}_{\text{NEW}}(x^n \rightarrow x^{n+1})$$

が成り立つ．

これまでの私たちの研究成果 [23], [24] では，接尾辞をタイプ I とタイプ II の 2 種類に分割することで， $\mathcal{A}_{\text{DEL}}(x^n \rightarrow x^{n+1})$ と $\mathcal{A}_{\text{NEW}}(x^n \rightarrow x^{n+1})$ の性質を示している．

ここで、タイプIの接尾辞とは、 x^n 上に1箇所だけ出現する接尾辞である．また、タイプIIの接尾辞とは、 x^n 上に複数箇所出現する接尾辞である．そこで、すべてのタイプIの接尾辞の集合を、

$$\mathcal{S}_I(x^n) = \{u \in \mathcal{S}(x^n) \mid u \notin \mathcal{D}(x^{n-1})\} \quad (5.3)$$

とおき、すべてのタイプIIの接尾辞の集合を、

$$\mathcal{S}_{II}(x^n) = \{u \in \mathcal{S}(x^n) \mid u \in \mathcal{D}(x^{n-1})\} \quad (5.4)$$

とおく．なお、 $\mathcal{S}_I(x^n)$ と $\mathcal{S}_{II}(x^n)$ は互いに素で、 $\mathcal{S}(x^n) = \mathcal{S}_I(x^n) \cup \mathcal{S}_{II}(x^n)$ であることを注記する．さらに、 $\mathcal{D}(x^{n-1})$ と $\mathcal{S}_I(x^n)$ は互いに素で、 $\mathcal{D}(x^n) = \mathcal{D}(x^{n-1}) \cup \mathcal{S}_I(x^n)$ であることを注記する．また、 $x^n \in \mathcal{S}_I(x^n)$ 、 $\lambda \in \mathcal{S}_{II}(x^n)$ である．

そして、 x^n の最も短いタイプIの接尾辞を

$$s_n = \operatorname{argmin} \{|v| \mid v \in \mathcal{S}_I(x^n)\} \quad (5.5)$$

と表記し、 s_n とタイプI・IIとの関係を次の命題5で示す．

命題5 任意の系列 x^n ($n \geq 1$) に対して、接尾辞 $u \in \mathcal{S}(x^n)$ が、

$$|u| \geq |s_n| \quad \text{ならば、} \quad u \in \mathcal{S}_I(x^n),$$

$$|u| < |s_n| \quad \text{ならば、} \quad u \in \mathcal{S}_{II}(x^n).$$

証明 最初に、 $|u| \geq |s_n|$ の場合を考える．このとき、 u と s_n は x^n の接尾辞なので、 $s_n \in \mathcal{S}(u)$ である．よって、 $s_n \in \mathcal{D}(u)$ となる．ところで、 s_n はタイプIの接尾辞なので、 $s_n \notin \mathcal{D}(x^{n-1})$ である．すなわち、 u 上に s_n が出現するから、 $u \notin \mathcal{D}(x^{n-1})$ となる．したがって、 $u \in \mathcal{S}_I(x^n)$ である．

次に、 $|u| < |s_n|$ の場合を考える．接尾辞 s_n は最も短いタイプIの接尾辞である．よって、 $u \notin \mathcal{S}_I(x^n)$ となる．したがって、 $u \in \mathcal{S}_{II}(x^n)$ である．

以上より、題意が成立する．

□

系5 任意の接尾辞 $u \in \mathcal{S}(x^n)$ が、 $u \in \mathcal{S}_I(x^n)$ 、 $\sigma(u) \in \mathcal{S}_{II}(x^n)$ ならば、 $u = s_n$ ．

次の命題Aと命題Bで、接尾辞のタイプと $\mathcal{A}_{\text{DEL}}(x^n \rightarrow x^{n+1})$ 、 $\mathcal{A}_{\text{NEW}}(x^n \rightarrow x^{n+1})$ との関係を示す．

命題A (定理6.5 [23])

任意の系列 x^n ($n \geq 1$) に対して、

$$\mathcal{A}_{\text{DEL}}(x^n \rightarrow x^{n+1}) = \{s_{n+1}\}.$$

命題 B (Proposition 3 [24])

任意の系列 x^n ($n \geq 1$) に対して, 系列 w が $w \in \mathcal{A}_{\text{NEW}}(x^n \rightarrow x^{n+1})$ ならば,

$$\pi(w) \in \mathcal{S}_I(x^{n+1}) \quad \text{または,} \quad \sigma(w) \in \mathcal{S}_I(x^{n+1}).$$

なお, 両方の命題とも本論文に適した形で表記を変えているが, 本質的に文献 [23][24] と同じ結果である.

5.5 新しく加わる極小禁止語の長さ

入力系列が x^n から x^{n+1} に変わったとき, 命題 B から, 反辞書 $\mathcal{A}(x^n)$ に加える極小禁止語の系列長は x^{n+1} のタイプ I の接尾辞の系列長に関係があることがわかる. そして, タイプ I の接尾辞の系列長は $|s_{n+1}|$ 以上である.

ここで, s_n と s_{n+1} の系列長の関係について, 次の補題 8 で示す.

補題 8 任意の系列 x^n ($n \geq 1$) に対して,

$$|s_{n+1}| \leq |s_n| + 1.$$

証明 命題 5 から, $\sigma(s_{n+1}) \in \mathcal{S}_{\text{II}}(x^{n+1})$ である. よって, $\sigma(s_{n+1}) \in \mathcal{D}(x^n)$. ここで, s_{n+1} の末尾は x_{n+1} だから, $u = \pi(s_{n+1})$ とおくと, $\sigma(u)x_{n+1} \in \mathcal{D}(x^n)$ となる. したがって, $\sigma(u) \in \mathcal{D}(x^{n-1})$ である. ここで, $\sigma(u)$ は x^n の接尾辞なので, $\sigma(u) \in \mathcal{S}_{\text{II}}(x^n)$ となる. したがって, $|\sigma(u)| < |s_n|$ である. すなわち, $|u| = |\sigma(u)| + 1 \leq |s_n|$ となる.

以上より, $|s_{n+1}| = |ux_{n+1}| = |u| + 1 \leq |s_n| + 1$ となり, 題意が成立する. 等号は, u が x^n の接尾辞なので, $u = s_n$ のとき, すなわち, $s_{n+1} = s_n x_{n+1}$ のときに限り成り立つ.

□

さらに, 新しく加わる極小禁止語の長さの上限と下限を, 定理 5 に示す.

定理 5 任意の系列 x^n ($n \geq 1$) に対して, 系列 w が $w \in \mathcal{A}_{\text{NEW}}(x^n \rightarrow x^{n+1})$ ならば,

$$|s_{n+1}| + 1 \leq |w| \leq |s_n| + 2.$$

証明 最初に, $|s_{n+1}| + 1 \leq |w|$ を証明する.

命題 B から, $\pi(w) \in S_I(x^{n+1})$ または, $\sigma(w) \in S_I(x^{n+1})$ である. すなわち, s_{n+1} は x^{n+1} のタイプ I の最短の接尾辞だから, $|\pi(w)| = |\sigma(w)| \geq |s_{n+1}|$ である. したがって, $|w| = |\pi(w)| + 1 = |\sigma(w)| + 1 \geq |s_{n+1}| + 1$ である.

次に, $|w| \leq |s_n| + 2$ を証明する.

命題 B から, $\pi(w) \in S_I(x^{n+1})$ の場合と, $\pi(w) \notin S_I(x^{n+1})$ の場合に分けて証明する. また, 表記を簡潔にするために $u = \pi(\sigma(w))$ とおく. ここで, w は x^{n+1} の極小禁止語なので, $\sigma(w) \in \mathcal{D}(x^{n+1})$ である. よって, $u \in \mathcal{D}(x^n)$ である.

まずは, $\pi(w) \in S_I(x^{n+1})$ の場合を考える. このとき, $u = \sigma(\pi(w)) \in S(x^{n+1})$ である. さらに, $u \in \mathcal{D}(x^n)$ から, $u \in S_{II}(x^{n+1})$ である. よって, 系 5 から, $\pi(w) = s_{n+1}$ である. このとき, 補題 8 から, $|\pi(w)| = |s_{n+1}| \leq |s_n| + 1$ である. したがって, $|w| = |\pi(w)| + 1 \leq |s_n| + 2$ となる.

今度は, $\pi(w) \notin S_I(x^{n+1})$ の場合を考える. このとき, 命題 B から, $\sigma(w) \in S_I(x^{n+1})$ である. ここで, もし $\pi(w) \in S(x^n)$ と仮定すると, $\sigma(w) \in S_I(x^{n+1})$ から, $w \in S(x^{n+1})$ となる. すなわち, $w \in \mathcal{D}(x^{n+1})$ である. これは, $w \in \mathcal{A}_{\text{NEW}}(x^n \rightarrow x^{n+1})$ に矛盾する. したがって, $\pi(w) \notin S(x^n)$ である. すなわち, $\pi(w) \in \mathcal{D}(x^{n-1})$ となる.

ところで, $\sigma(w) \in S_I(x^{n+1})$ から, $u = \pi(\sigma(w)) \in S(x^n)$ である. また, $\pi(w) \in \mathcal{D}(x^{n-1})$ から, $u = \sigma(\pi(w)) \in \mathcal{D}(x^{n-1})$ となる. よって, $u \in S_{II}(x^n)$ である. したがって, $|u| < |s_n|$ から, $|w| = |u| + 2 < |s_n| + 2$ である.

以上より, 題意が成立する. 等号は, $s_{n+1} = s_n x_{n+1}$ のときに限り成り立つ.

□

なお, 補題 8 と定理 5 の等号条件 $s_{n+1} = s_n x_{n+1}$ には s_{n+1} と x_{n+1} が使われている. そのため, 素朴な方法では x_{n+1} の入力後でないと等号条件を満たすかどうか分からない.

だが, 反辞書を用いることで, x_{n+1} の入力前に, どの記号が x_{n+1} になれば等号条件を満たすのかがわかることを, 次の命題 6 で示す.

命題 6 任意の系列 x^n ($n \geq 1$) に対して, $s_n x_{n+1} \in \mathcal{A}(x^n)$ ならば,

$$s_{n+1} = s_n x_{n+1}.$$

証明 系列 $s_n x_{n+1}$ は x^n の極小禁止語だから, $s_n x_{n+1} \notin \mathcal{D}(x^n)$ である. したがって, $s_n x_{n+1} \in S_I(x^{n+1})$ である. また, 同じく $s_n x_{n+1} \in \mathcal{A}(x^n)$ より, $\sigma(s_n) x_{n+1} \in \mathcal{D}(x^n)$ である. すなわち, $\sigma(s_n) x_{n+1} \in S_{II}(x^{n+1})$ となる.

よって, 系 5 より, 題意は成立する.

□

5.6 反辞書の変化と提案アルゴリズム

反辞書を動的に構築する手法は，これまでの私たちの研究成果 [23], [24] で提案している．しかし，従来の反辞書構築手法では，静的な手法でも動的な手法でも，入力記号系列から接尾辞トライや接尾辞木や接尾辞配列といった接尾辞のデータ構造を構築し，それから反辞書を構築している．

本論文では，この接尾辞のデータ構造の構築を省き，入力系列から直接反辞書オートマトンを構築する手法を提案する．

まず，命題 B により，新しく加わる極小禁止語は2つのタイプ $\mathcal{A}_{\text{NEW}}^{\text{I}}(x^n \rightarrow x^{n+1})$ ， $\mathcal{A}_{\text{NEW}}^{\text{II}}(x^n \rightarrow x^{n+1})$ に分けることができる．

$$\begin{aligned}\mathcal{A}_{\text{NEW}}^{\text{I}}(x^n \rightarrow x^{n+1}) &= \{w \in \mathcal{A}_{\text{NEW}}(x^n \rightarrow x^{n+1}) \mid \pi(w) \in S_{\text{I}}(x^{n+1})\}, \\ \mathcal{A}_{\text{NEW}}^{\text{II}}(x^n \rightarrow x^{n+1}) &= \{w \in \mathcal{A}_{\text{NEW}}(x^n \rightarrow x^{n+1}) \mid \sigma(w) \in S_{\text{I}}(x^{n+1})\},\end{aligned}$$

なお，両方のタイプに属する極小禁止語が存在する事もあるが， $\mathcal{A}_{\text{NEW}}(x^n \rightarrow x^{n+1}) = \mathcal{A}_{\text{NEW}}^{\text{I}}(x^n \rightarrow x^{n+1}) \cup \mathcal{A}_{\text{NEW}}^{\text{II}}(x^n \rightarrow x^{n+1})$ であることを注記する．

また， $\mathcal{A}_{\text{NEW}}^{\text{I}}(x^n \rightarrow x^{n+1})$ に属する極小禁止語を x^{n+1} のタイプ I の極小禁止語と， $\mathcal{A}_{\text{NEW}}^{\text{II}}(x^n \rightarrow x^{n+1})$ に属する極小禁止語を x^{n+1} のタイプ II の極小禁止語と呼ぶことにする．

さらに，定理 1 と系 1 の反辞書の別表現と組み合わせると，

$$\mathcal{A}_{\text{NEW}}^{\text{I}}(x^n \rightarrow x^{n+1}) = \{u\beta \mid \beta \in \mathcal{T}_{n+1}(u), u \in S_{\text{I}}(x^{n+1})\}, \quad (5.6)$$

$$\mathcal{A}_{\text{NEW}}^{\text{II}}(x^n \rightarrow x^{n+1}) = \{\alpha u \mid \alpha \in \mathcal{H}_{n+1}(u), u \in S_{\text{I}}(x^{n+1})\} \quad (5.7)$$

となる．

ここで，タイプ I の極小禁止語について，以下の補題 9 と定理 6 を示す．

補題 9 任意の系列 w に対して， $w \in \mathcal{A}_{\text{NEW}}^{\text{I}}(x^n \rightarrow x^{n+1})$ ならば，

$$\pi(w) = s_{n+1}.$$

証明 表記を簡潔にするために， $u = \pi(w)$ とおき，仮定法を用いて証明する．

もし， $u \neq s_{n+1}$ と仮定したとき， $u \in S_{\text{I}}(x^{n+1})$ と命題 5 から， $\sigma(u) \in S_{\text{I}}(x^{n+1})$ である．よって， $\mathcal{X}_{n+1}^{\text{T}}(u) = \mathcal{X}_{n+1}^{\text{T}}(\sigma(u)) = \phi$ となる．また， $u \neq \lambda$ したがって， $\mathcal{T}_{n+1}(u) = \mathcal{X}_{n+1}^{\text{T}}(\sigma(u)) \setminus \mathcal{X}_{n+1}^{\text{T}}(u) = \phi$ になる．これは， $w \in \mathcal{A}_{\text{NEW}}^{\text{I}}(x^n \rightarrow x^{n+1})$ に矛盾する．よって， $u = s_{n+1}$ である．

□

定理 6 任意の系列 x^n ($n \geq 1$) に対して ,

$$\mathcal{A}_{\text{NEW}}^{\text{I}}(x^n \rightarrow x^{n+1}) = \{s_{n+1}\beta \mid \beta \in \mathcal{T}_{n+1}(s_{n+1})\}.$$

証明 補題 9 と (5.6) から , 明らかである .

□

次に , タイプ II の極小禁止語について , 以下の定理 7 を示す .

定理 7 任意の系列 x^n ($n \geq 1$) に対して ,

$$\mathcal{A}_{\text{NEW}}^{\text{II}}(x^n \rightarrow x^{n+1}) = \{\alpha u \mid \alpha \in \mathcal{H}_{n+1}(u), u \in \mathcal{S}_{\text{I}}(x^{n+1}), |s_{n+1}| \leq |u| \leq |s_n| + 1\}.$$

証明 定理 5 と (5.7) から , 明らかである .

□

ところで , 系列 x^n の最短のタイプ I の接尾辞 s_n は , x^n の反辞書オートマトンのノード $\mathcal{N}(x^n)$ のラベル列と一致する . また , 命題 A から , s_{n+1} は MFW ノード $\mathcal{N}(x^n x_{n+1})$ のラベル列である . この MFW ノードを p とおく . 定理 6 を用いると , p の suffix link 先ノード $\mathcal{S}(p)$ から記号 β で遷移した先が中間ノードのとき , $L(p)\beta$ は x^{n+1} のタイプ I の極小禁止語である .

一方 , x^{n+1} のタイプ II の極小禁止語 $w \in \mathcal{A}_{\text{NEW}}^{\text{II}}(x^n \rightarrow x^{n+1})$ は , 定理 7 から , $\sigma(w) \in \mathcal{S}_{\text{I}}(x^{n+1})$, $|s_{n+1}| \leq |\sigma(w)| \leq |s_n| + 1$ である . さらに , 命題 5 から , $\sigma(w)$ の末尾に s_{n+1} が出現する . すなわち , w を x^n の反辞書オートマトンに入力すると , p に到達する . したがって , ルートから p へ至る全てのパスから w を見つけばよい .

以上のことから , 新しい記号 x_{n+1} が入力されたときに反辞書オートマトンを変化させるアルゴリズム change_AAM を以下に提案する .

このアルゴリズムは , 既に x^n の反辞書オートマトンが構築済みであるものとする . また , $x^0 = \lambda$ の反辞書オートマトンは , 初期化作業により与えられるものとする .

Algorithm change_AAM

```

begin 1
   $p \leftarrow \mathcal{N}(s_n x_{n+1})$  2
   $p$  からタイプ I の MFW ノードを伸ばす . 3
  for (  $\forall w : p$  に至るパスから生成した系列 ) do 4
    if (  $|s_{n+1}| \leq |\sigma(w)| \leq |s_n| + 1$  ) 5
```

w をラベル列に持つ新しい MFW ノードを追加する .	6
end.	7

第6章 おわりに

反辞書は，接尾辞木や接尾辞配列と同様に，入力系列の特徴が反映されている．そして，反辞書オートマトンにより，計算機への負荷が少なく高速に動作するという利点がある．データ圧縮などでも，優れた圧縮性能を持っている．しかし，反辞書構築による計算機への負荷がとて大きいため，扱いづらい面もある．

そこで，本論文では，効率的に反辞書を構築する手法を検討した．

得られた結果は，以下のとおりである．

第4章では，接尾辞配列を用いた静的な反辞書構築手法を提案した．検討の過程では，反辞書の別表現を新たに提案し，極小禁止語を，先頭記号と後続する系列に分けて考えることにより，あと1記号先頭に加えることで極小禁止語になる系列が，入力系列の接尾辞配列ならびにL配列によって特徴づけられること，すなわち，先頭に1記号加えれば極小禁止語になる系列は，接尾辞配列上の接尾辞の先頭部分に存在し，その長さはL配列で求めることができることを明らかにした．これにより，線形時間での構築が可能となった．実験結果では，接尾辞木を用いた構築手法と比べて，記憶量と計算量を改善することができた．

第5章では，動的に反辞書オートマトンを構築する手法を提案した．従来の反辞書構築手法では，静的手法でも動的手法でも，入力記号系列から接尾辞トライや接尾辞木や接尾辞配列などを構築し，それから反辞書を構築する．そこで，本論文では，入力系列の末尾に新たな記号が加わるたびに，逐次的に反辞書オートマトンを変化させ，入力系列から直接反辞書を構築する手法を提案する．

そのために，入力系列の末尾に新しい記号が追加されることで，反辞書にどのような変化が起きるのかを明らかにした．これにより，接尾辞のデータ構造を経由しないで反辞書オートマトンを構築することが可能となった．

今後の研究課題としては，この成果を活用した応用技術を開発することである．

謝辞

本研究を御指導して下さった電気通信大学情報システム学研究科 森田 啓義教授には、研究が遅々として進まないだけでなく、執筆も遅い私へ親身に助言と御指導をして戴きました。そして、研究方法と研究者の心構えについても多くのことを学びました。このことに深く感謝しお礼申し上げます。

また、本研究を通してとても貴重な助言を戴きました長野県工科短期大学 校 太田 隆博 准教授にも、多くの迷惑を掛けながらも、親身に助言を戴きました。心より厚くお礼申し上げます。

そして、研究のきっかけを与えてくださった職業能力開発総合大学校の松嶋 智子 准教授、研究しやすい環境づくりに影から支えてくださった高知職業能力開発短期大学校の松中 孝二 校長をはじめとする職場の同僚に感謝いたします。

最後に、何年間も高知と東京の往復をすることを支えてくれ、ときには挫折しそうなときに叱咤激励してくれた家族に深く深く感謝いたします。

ありがとうございました。

関連論文リスト

学術雑誌論文（査読あり）

- [1] 深江 裕忠, 太田 隆博, 森田 啓義 : ”配列を用いた反辞書の線形構築法”, 電子情報通信学会和文論文誌 A Vol.J97-A,No.1, pp.33–44 (2014). 【第4章と関係】

国際会議発表（査読あり）

- [2] H. Fukae, T. Ota and H. Morita : ”On fast and memory-efficient construction of an antidictionary array”, Proc. of 2012 IEEE Int’l Symp. on Information Theory (ISIT2012), Cambridge, pp.1092–1096 (2012). 【第4章と関係】

口頭発表

- [3] 深江 裕忠, 太田 隆博, 森田 啓義 : ”接尾辞配列による反辞書構築法”, 第33回情報理論とその応用シンポジウム (SITA2010), 松代, pp.829–834 (2010). 【第4章と関係】
- [4] 深江 裕忠, 森田 啓義 : ”Suffix array を利用した動的な反辞書生成”, 第30回情報理論とその応用シンポジウム (SITA2007), 賢島, pp.384–387 (2007). 【第5章と関係】

参考文献

- [5] M. Crochemore, F. Mignosi, A. Restivo and S. Salemi : "Data compression using antidictionaries", Proc. IEEE, 88, 11, pp.1756–1768 (2000).
- [6] 西新, 森田, 太田 : "反辞書木を用いた分岐予測手法", 電子情報通信学会論文誌 A Vol.J93-A No.3, pp.204–206 (2010).
- [7] D. Sun, H. Morita and M. Nishiara : "On construction of reversible variable-length codes including synchronization markers as codewords", Proc. of 2006 Int'l Symp. on Information Theory and its Applications (ISITA2006), pp.12–16 (2006).
- [8] 太田, 戸松, 木村, 森田 : "反辞書符号化法を用いた不整脈検出", 信学技報 MBE2010-7, pp.35–40 (2010).
- [9] T. Ota and H. Morita : "On the construction of an antidictionary with linear complexity using the suffix tree", IEICE Trans. on Fundamentals, E90-A, 11, pp.2533–2539 (2007).
- [10] U. Manber and G. Myers : "Suffix arrays: A new method for on-line string searches", SIAM Journal of Computing, 22(5), pp.935–948 (1993).
- [11] 笠井 : "部分語計数問題の接尾辞配列を用いた高速アルゴリズム", Master's thesis, 九州大学大学院システム情報科学研究科 (1999).
- [12] T. Kasai, G. Lee, H. Arimura, S. Arikawa and K. Park : "Linear-Time Longest-Common-Prefix Computation in Suffix Arrays and Its Applications", Proc. of the 12th Annual Symposium on Combinatorial Pattern Matching (CPM 2001), Lecture Notes in Computer Science, 2089, pp.181–192 (2001).
- [13] A. Aho ほか著, 大野義夫 訳 : "データ構造とアルゴリズム", 培風館 (1987).
- [14] D. Gusfield : "Algorithms on strings, trees, and sequences", Cambridge Univ.Press (1997).
- [15] E. Ukkonen : "On-line construction of suffix-trees", Algorithmica, 14, pp.249–260 (1995).

- [16] G. Nong, S. Zhang and W. H. Chan : "Linear Suffix Array Construction by Almost Pure Induced-Sorting", Proc. of 19th IEEE Data Compression Conference (2009).
- [17] J. Kärkkäinen and P. Sanders : "Simple Linear Work Suffix Array Construction", Proc. of 30th International Conference on Automata, Languages and Programming, pp.943-955 (2003).
- [18] J. Kärkkäinen, P. Sanders and S. Burkhardt : "Linear Work Suffix Array Construction", J.ACM 53 (6), pp.918–936 (2006).
- [19] D. K. Kim, J. S. Sim, H. Park and K. Park : "Constructing Suffix Arrays in Linear Time", J.Discrete Algorithms 3(2-4), pp.126–142 (2005).
- [20] P. Ko and S. Aluru : "Space Efficient Linear Time Construction of Suffix Arrays", J.Discrete Algorithms 3(2-4), pp.143–156 (2005).
- [21] R. Sedgwick : "Algorithms in C", Addison-Wesley (1990).
- [22] T. Asano, S. Bereg and D. Kirkpatrick : "Finding Nearest Larger Neighbors – A Case Study in Algorithm Design and Analysis –", Lecture Notes in Computer Science, 5760, pp.249–260 (2009).
- [23] 深江 : "接尾辞配列を利用した動的反辞書生成に関する研究", Master's thesis, 電気通信大学大学院情報システム学研究科 (2008).
- [24] T. Ota, H. Morita and H. Fukae : "On the Dynamic Construction of an Anti-dictionary with Linear Complexity", Proc. of 2011 IEEE Int'l Symp. on Information Theory (ISIT2011), pp.2422–2426 (2011).